

BEST AVAILABLE COPY

5章

ファイルシステム

オペレーティング・システムにおいて、最も頻りに利用されるのはファイルシステムである。プログラムでは、ファイルの読み書きが行われ、またユーザーはファイルの存在とその属性を常に意識している。オペレーティング・システムが使いやすく、そして便利なものであるかどうかは、ファイルシステムのインターフェイス、構造、信頼性に大きく依存している。

本章では、ユーザー側から見たファイルシステムのさまざまな形態、その実装方法、不正使用に対するファイルの保護、MINIX ファイルシステムの設計と実装方法などを解説する。また分散オペレーティング・システムにおいて見られる設計上の問題点についても解説していく。

5.1 ユーザーから見たファイルシステム

ユーザーの立場において、ファイルシステムで最も重要なものは、その形態である。すなわちファイルが構成しているものが何か、その命名法、保護法、許可されるファイル操作、などが最大の関心事である。ファイルシステムの詳細、例えば空き領域を管理するために連結されたリストやビットマップをどの様に使用しているか、また論理ブロック内のセクタ数など、細かいことはそれほど重要ではない。ただし、ファイルシステムの設計者にとっては、この様な詳細が非常に重要である。以後の節では、ユーザーインターフェイスに関するいくつかの事項を説明する。章の後半ではファイルシステムの構築法に関して考察する。

297

4章 メモリ管理

16. 最初のタイムシェアリング・マシンの1つ、PDP-1は18ビットワードを4K持っていた。メモリ内に常時1プロセスを保存していた。スケジューラが別のプロセスを実行しようとする時は、メモリ内のプロセスはページ・ドラムに書き込まれ、そのドラムには1回あたり18ビットワードが4K記録できた。しかもドラムはワード0からだけでなく、任意のワードから読み書きすることができた。このドラムがどうして使用されることになったのかを考えよ。

17. あるコンピュータは各プロセスに対して4096バイトページに分割された65,536バイトのアドレス空間を割り当てている。特定のプログラムが32,768バイトのテキストサイズ、16,386バイトのデータサイズ、そして15,870バイトのスタックサイズを持っているとしよう。このプログラムはアドレス空間に適合するか。ページサイズが512バイトならどうか。1つのページは2つの異なるセグメントを含まないことに注意せよ。

18. ページフォルトとページフォルトの間に実行された命令数は、プログラムに割り当てられたページフレームの数に正比例していることがわかっている。使用可能なメモリが倍になれば、ページフォルトと間の平均間隔も倍になる。通常の命令は1マイクロ秒しかかからないが、ページフォルトが発生すると2001マイクロ秒かかってしまう。プログラムの実行に60秒かかり、その間に15,000のページフォルトが発生するとして、この時2倍のメモリ量を与えた場合、実行時間を求めよ。

19. MINIX のメモリ管理手法では、chmemの様なプログラムがなぜ必要なのか説明せよ。

20. MINIX プロセスがSIGNAL システムコールを発行し、シグナルを捕えられようとした時、呼び出される関数のアドレスはメモリマネージャではなく、プロセス自身のテーブルに保存されることになる。この様な技法を用いることによる利点をあげよ。

21. 関プロセスが作機共題になるまで待たずに、ゾンビ状態になったプロセスのメモリを即座に解放できるようにMINIXを修正せよ。

22. 現在のところMINIXでは、EXECシステムコールが発行されると、メモリマネージャが新しいメモリノージを保存するのに十分な大きさのホールが存在しているかを調べる。それだけの大きさのホールが存在しなければ、呼出しは拒絶される。より高度なアルゴリズムでは、現在使用中のメモリノージが解放された時、十分な大きさのホールが存在するかを調べることである。このアルゴリズムを実現せよ。

23. EXECシステムコール実行時に、MINIXはファイルシステムにセグメント全体を一度に読み取らせるようなトリックを行う。同じような方法でコアダンプを書き込むためのトリック技法を説明し、実装せよ。

24. MINIXを修正し、スワッピングを可能にせよ。

296

5.1.1 ファイルに関する基礎

コンピュータに情報を保存するための最速な方法としては、各プロセスに非常に大きなセグメントを多数、例えば2²⁰の長さを持つセグメントを2²⁰個与えることである。誰かが最初にコンピュータを使い始めるとき、シミュレーションのセグメント(ここではエディタ、コンパイラ、その他ユーティリティなどが含まれる)のためのアドレス空間が与えられる。時間が経過するにつれ、ユーザーの入力したテキスト、プログラム、その他の情報はすべてシミュレーションのアドレス空間内に蓄積されていく。あるセグメントにはディレクトリが格納されており、階層構造のファイル管理(UNIXの様な)を行っている。

このアドレス空間は、すべてのプロセスによって継承され、必要があればセグメントを追加し、不要になればアドレス空間はシミュレーションに返される。これにより各情報のアクセスは、メモリ(仮想)に対する読み書きだけの問題となる。

MULTICSではこの問題に真剣に取り組んだが、実際にはこの様な情報の格納方法はまだ実現されていない。その理由の1つには、現在のアドレス空間が小さすぎ、すべてのコンピュタが仮想アドレスを持っていておけないということがある。また、プロセスがクラッシュした場合、そのアドレス空間が失われ、アドレス空間を長期間(例えば月や年単位で)保存できなくなる、という理由の1つである。それでも技術が進歩すれば、ユーザーのプロセスのアドレス空間にユーザーの情報をすべて保存しておくという概念は、再度検討されるだろう。

その様な環境が整うまで、大半のオペレーティングシステムでは、これは適当な方法を用いて情報を保存するだろう。それは、ユーザーが名前付けできるファイル(files)を用いるという方法である。ファイルに、プログラムやデータ、その他の必要なものを保存する。ファイルは、どのプロセスのアドレス空間の一部でもない、オペレーティングシステムは、ファイルの作成、削除、読み書きを行う特殊な操作(システムコール)を提供しており、またファイルの管理している。

図5.1は一般的な3つのファイルの構成を示したものである。最初の方法は、単なるバイナリ列である。UNIXのファイルは、この方法で構成されている。2番目の方法は、固定サイズ長の連続したレコードである。任意のレコードを読み書きできるが、ファイルの途中にレコードを挿入する、あるいは削除するなどは行えない。CP/Mは、この方法を採用している。

3番目の方法は、ツリー構造のファイルシステムである。個々のブロックがn個のキーを持つレコードを持っている。レコードはキーを用いて検索することができ、新しいレコードはツリー内の任意の場所に挿入することができ、すでに満杯となっているブロックに対し、新たにレコードが加えられた場合、ブロックは2つに分断され、両方ともアルファベット順に正しく並べられたツリーに加えられる。この方法はメインフレームにおいてよく採用されており、ISAM(Indexed Sequential Access Method)と呼ばれている。

オペレーティングシステムでは、デバイスに依存しない(device independence)ことが重要である。すなわちファイルまたは装置がどこに存在しているかと、同じアクセスを行えるように

5.1 ユーザーから見たファイルシステム

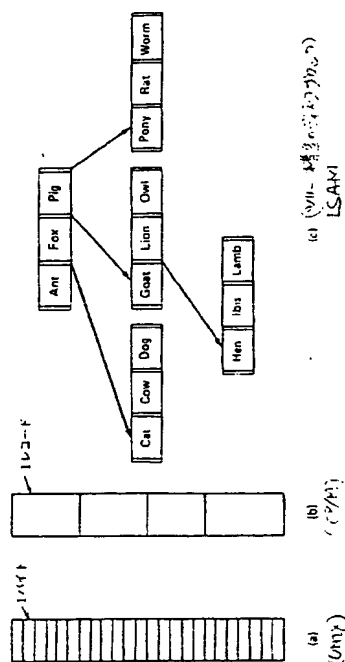


図 5.1 3種類のファイル

(a) バイト列 (b) レコード列 (c) ツリー

することである。例えば、入力ファイルを読み取り、それをソートし、結果を出力ファイルを書き込むプログラムでは、フロービディスタンスや、ディスタンスなど、どのデバイスからでもそれを意識することなく入力が行え、またファイル、端末、あるいはプリンタなど、どのデバイスに対しても、デバイスを意識することなく出力が行えるということである。

あるオペレーティングシステムでは、他のオペレーティングシステムと比較して、装置に対する依存性が低いものもある。例えば UNIX ではファイルシステム(例えばディスク)をファイルツリー内の任意の場所にマウントできるようにしており、すべてのファイルを、どの装置上にあるかを意識することなく、パス名を用いてアクセスできる。MS-DOS では、ユーザーがどの装置上にファイルが存在しているかを明示的に指定しなければならない(ただし、デフォルトの装置となっているものだけは必要ない)。したがって、デフォルトの装置がドライブCであれば、ドライブA上にあるプログラムをドライブBに置かれた入出力ファイルを用いて実行するために、以下の様な入力を行わなければならない。

A: program < B: input > B: output

大半のオペレーティングシステムでは、多くの異なるファイル型を持っている。例えば UNIX には、通常のファイル、ディレクトリ、特殊ファイル(ブロックおよびキャラクタ)がある。通常のファイルはユーザーデータを含んでいる。ディレクトリにはファイルに記号名(すなわち ASCII 文字列)を付けるために必要な情報が格納されている。ブロックおよびキャラクタ型特殊ファイルは、それぞれディスク装置および端末装置などを表すために用いられる。これにより、例えばユーザーは以下の様な入力を行うことによって、ファイルabcを端末(キャラクタ型特殊ファイル/dev/tty)にコピーすることができる。

5章 ファイルシステム

cp abc /dev/tty

UNIXのあるバージョンでは、名前付きパイプ(named pipe)と呼ばれる技法も用いられている。名前付きパイプは、2つのプロセスにおいてオープンし、そのプロセス同士で通信を行うことができる疑似ファイルである。ユーザーからすれば、これらのファイル型の違いは些細なものである(例えば端末に対してシークを行うことができないなど)。

大半のシステムでは、通常のファイルがさらにその用途に応じて異なる型に分類されている。個々の型は、ファイル名の最後に付けられるファイル拡張子(file extensions)によって識別される。以下にその一例を示す。

```
FILE.PAS Pascal ソースプログラム
FILE.FTN FORTRAN ソースプログラム
FILE.BAS BASIC ソースプログラム
FILE.OBJ オブジェクトファイル
FILE.BIN 実行可能バイナリファイル
FILE.LIB リンカの使用済みOBJ ファイルのライブラリ
FILE.TXT フォーマットされたテキストファイル
FILE.DAT データファイル
```

システムによっては拡張子を単なる慣習として使用しており、オペレーティング・システム自身はそれを特定の目的に使用していない。また、オペレーティング・システムが命名に関する厳しい制限を設けているシステムもある。例えば、BIN で終了しない限りファイルを実行しないことになっている、という例もある。

TOPS-20 システムでは、実行されるファイルの作成時間の管理まで行っている。そしてソース・ファイルを探し、バイナリが作成されてから後ソースが修正されたかを調べている。修正されれば、自動的にソースをコンパイルし直す。UNIX システムにおける make プログラムに該当するものが、シェルスクリプトに内蔵されているのである。拡張子が必ず必要であるため、オペレーティング・システムは、どのバイナリプログラムがどのソースから作成されているのかを識別できる。この様に、システムに対して強い型付けを行うと、ユーザーがシステム設計者の予測しなかったことを行おうとするために問題が発生してしまう。例えば出力結果として、DAT(データファイル)型のファイルを作成するプログラムを考えてみる。ユーザーがPAS ファイルを読み取り、それを変換し(例えば標準のインデント形式に変換する)、変換したファイル出力へ書き込む場合、出力ファイルの型はDAT となる。これを Pascal コンパイラを用いてコンパイルすると、拡張子が誤っているため、システムはコンパイルを行わない。FILE.DAT を FILE.PAS へコピーしようとしても、システムはコピーを拒否する(ユーザーの誤りを防ぐため)。

5.1 ユーザーから見たファイルシステム

オペレーティング・システムの判断によってこの様な保護が行われることは、初心者にとってはありがたいものであるが、上級者にとってはうとましいものである。

ファイルに対して行える操作は、オペレーティング・システムによって異なっているが、バイト(あるいはレコード)の連続的な読み書きはどのシステムでも行える。ランダムアクセス(random access)は、効果的であると思われた場合に提供されている(通常のファイルに対しては、READ が、ラインプリンタの様な特殊ファイルでは不可)。あるオペレーティング・システムでは、READ システムコールを使用する際に読み書きするレコードの番号(あるいはキー)を指定し、ランダムアクセスが行われている。また、ファイル内における「現在のファイル位置」を設定するシステムコールが提供されており、それにより以後の READ または WRITE コールが新しい位置に対する操作を行う。というシステムもある(例えば UNIX など)。レコードを用いているファイルでは、レコードの挿入と削除の操作が行えるようになっている。

5.1.2 ディレクトリ

ファイルを記録しておくために、ファイルシステムでは通常ディレクトリ(directory)を提供しており、大半のシステムではこのディレクトリ自身もファイルの一種となっている。ディレクトリには通常いくつかのエントリがファイルごとに含まれている。これを示したのが図 5.2 である。最も簡単な方法では、システムがすべてのユーザーの全ファイルを持つ1つのディレクトリを管理している。多くのユーザーが存在する場合、それぞれのユーザーが同じファイル名を使う(例えば mail や games など)。衝突や、混乱が生じ、システムが稼働できなくなる。このようなシステムは最も初歩的なマイクログコンピュータ用オペレーティング・システムにしか用いられていない。

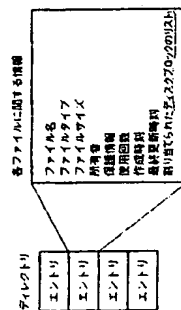


図 5.2 ディレクトリ・エントリ

ディレクトリにはいくつかのファイルに関するエントリを、複数のファイル分持たせることができる。エントリにはファイルに関する情報や、その他の情報を格納しておく他の構造体を持つポインタが含まれている。

5.2 ファイルシステム

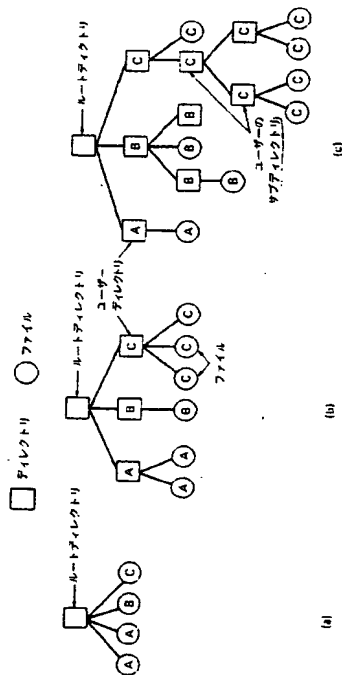


図 5.3 3つのファイルシステム設計

- (a) すべてのユーザーが1つのディレクトリを共有している
- (b) 1ユーザーあたり1ディレクトリ
- (c) ユーザーごとの、図中の文字はファイルやディレクトリの所有者を示す

1つのディレクトリにすべてのファイルを格納するという考えに改良を加え、1ユーザーあたり1つのディレクトリを持たせることもできる(図5.3(b)参照)。この方法はユーザー間でのファイル名の衝突を防ぐことにはなったが、多くのファイルを用いるユーザーにはあまり効果的ではなかった。論理的な方法で自分のファイルをグループ化しようとするユーザーも少なくなかった。例えば大学教授の場合、ある講義のために作成している資料のファイル、またそれをさらに本にまとめるファイル、別な講義に使用する学生用のプログラム、自分自身用に作成している高度なコンパイル作成プログラム、承認済みの提案を格納しているファイル、メール送受信用のファイル、講義録、書きかけの論文、ゲームなど、多種多様なファイルを持っており、これらを整理するための様々なファイルグループが必要となる。これらのファイルを論理的にまとめるための方法が必要である。

ここで必要となるのは、階層化(すなわち「ディレクトリ・ツリー」)である。このアプローチにより、ユーザーはファイルが自然な方法でグループ分けできるように、必要なだけディレクトリを持つことができるようになる。これを示したのが図5.3(c)である。

ファイルシステムがディレクトリのツリーとして構成されると、ファイル名を指定する方法を決めておかねばならない。一般的には次の2つの方法が使用されている。最初の方法では、各ファイルに、ルートディレクトリからそのファイルまでのパスを表す絶対パス名(absolute path name)を付ける。例えばパス名/user/ast/mailboxは、ルートディレクトリにサブディレクトリuserが存在し、その中にさらにastというサブディレクトリが存在し、そこにファイルmailboxが含まれていることを示すものとなる。絶対パス名は、必ずルートディレクトリから始まり、一意な名前となる。

5.2 ファイルシステムの設計

もう1つの方法は相対パス名(relative path name)である。これはワーキングディレクトリ(working directory)の概念と共に使用する。ワーキングディレクトリはカレントディレクトリ(current directory)とも呼ばれている。ユーザーは1つのディレクトリをカレント・ワーキングディレクトリとして定義する。ルートディレクトリから始まらないパス名は、ワーキングディレクトリから相対的に解釈される。UNIXではスラッシュで始まるパス名はすべて絶対パス名であり、その他はすべてユーザーのカレント・ワーキングディレクトリからの相対パス名となる。ワーキングディレクトリが/user/astであれば、絶対パス名が/user/ast/mailboxのファイルは、単にmailboxと指定するだけでよい。

5.2 ファイルシステムの設計

それではファイルシステムをユーザーではなく、システム設計者の側から考えてみることにする。ユーザーはファイルの命名法、認められている操作内容、ディレクトリ・ツリーの形態、その他インターフェース関連の事項に関心があつた。設計者はディスク空間の管理法、ファイルの保存法、効率よく、信頼性の高いシステム構築法などが興味深い。以降の項では、問題点とそれらを解決するために必要とされるトレードオフについて説明する。

5.2.1 ディスク空間の管理

ファイルは通常ディスクに保存されている。したがってディスクの管理はファイルシステムを設計する者にとって最大の関心事となる。nバイトのファイルを保存するためには、2つの方法がある。連続するnバイト上のディスク空間を割り当てる方法と、ファイルをいくつかの断片化し、断片化している必要はないか?ブロックに分割する方法である。同じような問題が断片化をセグメンテーションとページングの間のメモリ管理システムにも存在している。

ファイルを連続するバイトとして保存することは明らかに問題である。ファイルが成長した場合には、ディスク上の別な位置に移動しなければならぬ可能性があるからである。同じ問題はメモリ内のセグメントにも該当する。ただし、ファイルをディスク上で移動するのに比べ、メモリ内におけるセグメントの移動は高速である。

したがってほとんどのファイルシステムでは、ファイルを固定サイズのブロックに分割し、かつブロックが断片化している必要がない、という方法をとっている。ファイルを固定サイズのブロックに分割し保存する場合には、ブロックの大きさの問題となる。ディスクの構造から考えて、セクタ、トラック、シリンドラなどが割当て単位として考えられる。ページングシステムではページサイズの利用も考えられる。

シリンドラの様に大きな割当て単位を用いしまうと、1バイトのファイルでもシリンドラ全体を占有してしまうことになる。UNIXにおける平均ファイルサイズは約1Kであることが研究の結

5 5.2 ファイルシステム

果(Mullender Tanenbaum, 1984) 判明している。したがって各ファイルに32Kのシリンダを割り当てるとは、31/32 すなわち97%のディスク空間を無駄にすることとなる。逆に小さな割当て単位を用いると、各ファイルが多くのブロックから構成されることになる。そうすると、各ブロックの読み出しは、シークと回転によるロスを生じるため、多くの小さなブロックから構成されるファイルを読み出す場合の速度が低下する。

一例として、トラックあたり32,768バイトのディスクを想定してみる。回転時間は16.67ミリ秒、平均シーク時間は30ミリ秒である。k バイトのランダムブロックを読み出すための時間をミリ秒で表すと、以下の様にシーク、回転による遅延、転送時間の合計となる。

$$30 + 8.3 + \frac{(N/32768) \times 16.67}{(k/2)}$$

図5.4の表例は、この様なディスクのデータ転送効率をブロックサイズの関数として示したものである。すべてのファイルが1Kであるとは仮定した場合(測定平均サイズ)。図5.4の点線はディスク空間の効率を表す。空間を効果的に使用すると(ブロックサイズ<2K)、データ転送効率が落ちるという矛盾が生じることになる。空間的効率と、時間的効率は互角に立する関係にある。

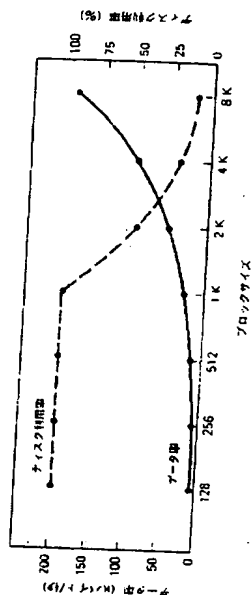


図5.4 ディスク空間利用率
実線(左側)はディスクのデータ転送効率を示したものである。点線(右側)はディスク空間の効率を示したものである。すべてのファイルは1Kとする。

通常の解決策は、512, 1K, または2Kのブロックサイズを選択するという方法である。1Kのブロックサイズと512バイトのセクタサイズを持つディスクで選択した場合、ファイルシステムは必ず2つの連続するセクタを読み書きし、それを1つの分割不可能な単位として扱うことになる。

ブロックサイズを決定した後は、空きブロックの管理方法を決めなくてはならない。一般的には図5.5に示す2つの方法が用いられている。最初の方法はそれ自身が持ちうる限りの空きブロックの番号を格納しているディスクブロックの連結されたリストを用いるものである。ブロックが1Kであり、ディスクブロック番号に16ビット使用している場合、空きリスト上の各ブロック

5.2 ファイルシステムの設計

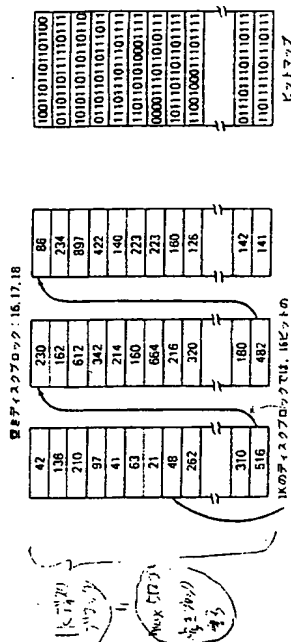


図5.5 空き領域の管理
(a) 連結されたリスト上に空きリストを持つ (b) ビットマップ

図5.5 空き領域の管理
(a) 連結されたリスト上に空きリストを持つ (b) ビットマップ

クには、512 個の空きブロック番号が格納されることになる。20M のディスクでは 20,000 のディスクブロック番号すべてを保存するために最大 40 ブロックの空きリストが必要となる。

空き領域管理において、ビットマップを用いる方法もある。n ブロックを持つディスク空間では、n ビットのビットマップが必要となる。マップにおいて、空きブロックは 1、また割り当てられたブロックは 0 で示される(またはその逆)。20M のディスクでは、マップに対して 20K ビット、すなわち 3 ブロックしか必要としない。ビットマップは 1 ブロックあたり 1 ビットを使用するだけでよい。連結されたリストの場合、1 ブロックあたり 16 ビットが必要であったが、これと比較すると、かなり少ない空間で済むわけである。

ビットマップを保存できるだけの主記憶があれば、この方法が適している。ただし、メモリ内において、空きディスクブロック・リストを格納する領域が 1 ブロックしか確保できず、またディスクもほぼ満杯である場合、連結されたリストを使っただろう。メモリ内に 1 ブロック分のビットマップを置く場合、もし空きブロックがそのビットマップ内で検出できなければ、ディスクをアクセスし、残りのビットマップを読み出さなくてはならない。連結されたリストであれば、新しいブロックがメモリにロードされると、次のブロックを読み取るまでに、512 のディスクブロックを割り当てることができる。

5.2.2 ファイルの保存

ファイルが連結したブロックから構成されている場合、ファイルシステムは各ファイルを構成するブロックを何らかの方法で管理しなくてはならない。ブロックを連続的に保存するという最も簡単な方法は、ファイルが成長していくため現実的ではない。事実、ファイルをブロックに分けざるを得なかったのもこの問題に起因している。

5.2 ファイルシステム

ファイルのブロックを連結されたリストとして保存することは可能である。1024バイトのディスクブロックに、1022バイトのデータと、次のブロックを示すポインタを含める方法である。これには2つの問題がある。まずブロック内のデータバイト数が2の階乗ではなくなくなってしまうという問題が生じる。さらに深刻であるのは、ランダムアクセスの実行において多くの無駄が発生することである。プログラムがファイル内の32768/1022、すなわち33ブロックをシークし、読み始めようとしている時、オペレーティングシステムは、32768/1022、すなわち33ブロックのシークし、必要なデータを抽出しなければならない。シークを行うために33ブロックもの読取りを行うことは、効率的ではない。

しかしファイルの連結されたリストとして表現する際の問題は、ポインタをメモリ内に保存することによって解決できる。図5.6はMS-DOSに用いられている割当て手法である。この例では、ブロック6、8、4、2が割り当てられているファイルA、ブロック5、9、12が割り当てられているファイルB、そしてブロック10、3、13が割り当てられているファイルCを示している。

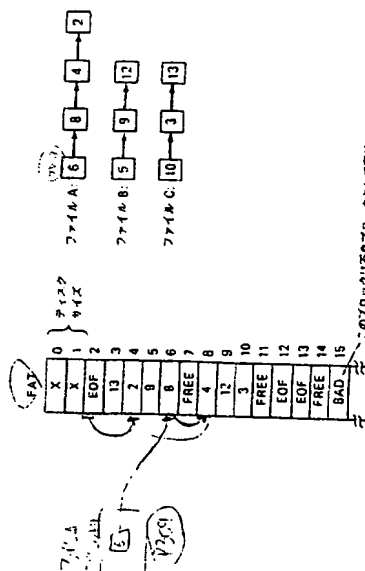


図 5.6 MS-DOS で用いられている連結されたリストの割当て手法
エントリ0と1はディスクサイズの指定に用いられる。EOF はファイルの終端、FREE は空きブロックである。

各ディスクに用意されるテーブルを、ファイル割当てテーブル (FAT: File Allocation Table) と呼ぶ。1つのエントリが1ディスクブロックに対応している。各ファイルのディレクトリには、ファイルの先頭ブロックのブロック番号が格納されている。そのブロック番号に対応する FAT のスロットには、次のブロックのブロック番号が格納されている。ファイルAはブロック6から始まり、FATのエントリ6はファイルAの次のブロック番号である8を格納している。FATエントリ8は次のブロック番号である4を格納している。エントリ4はエントリ2を、そしてエントリ2にはファイル終端 (End of File) の印が付けられる。

5.2 ファイルシステムの設計

この手法は、元来 1K のブロックサイズを用いた 320K のフロッピーディスク (MS-DOS 標準) に用い提供されたものであった。ブロック番号は 12 ビットであるため、320 エントリの FAT は 480 バイトを要することになり、512 バイトのセクタにちょうど納まる。IBM が DOS 2.0 から始めた、フロッピーディスクを 360K バイトとしてフォーマットするという方法では、FAT が 540 バイトに増加した。この結果、1つのセクタには格納できなくなり、ディスクのレイアウトを変更してより大きな (2 セクタ) FAT を作成しなければならなかった。4096 よりブロック数の多いハードディスクが出現すると、12 ビットのブロック番号は不適当となり、FAT を再度変更しなければならなくなった。先見の明があれば、変更を繰り返す必要もなかっただろう。

また大容量のディスクにおいては、この手法を採用することが少なくなっている。1K のディスクブロックを 64,000 個含む 64M のディスクを仮定してみる (70M のディスクであれば 16 ビット以上のブロックサイズが必要となり、さらに厄介なことになる)。FAT は 64K 分のエントリ (各 2 バイト) を持ち、128K を占有することになる。これをすべてメモリ内に常駐させると、かなりの量のメモリが消費されてしまう。しかしディスクに保存すると、ファイル内の 32K バイト H に対するシークを行う場合、ブロックの連鎖をたどるために 1 回以上 33 回以下のディスク読取りが必要となる。

FAT の基本的な問題は、ディスク上のすべてのファイルに対するポインタが、1つのテーブル内にランダムに散在しているという点である。したがってたとえ 1つのファイルをオープンする場合でも、FAT 全体が必要になる。よりよい方法として、各ファイルのディレクトリのブロックリストを別々の場所に保存しておく方法がある。UNIX はこの様な方法を採用している。

UNIX では各ファイルに対応して、インデックス (inode) と呼ぶ小さなテーブルが用意されている。図 5.7 はこれを示したものである。このテーブルには「カウント情報」と、保護情報が含まれている。

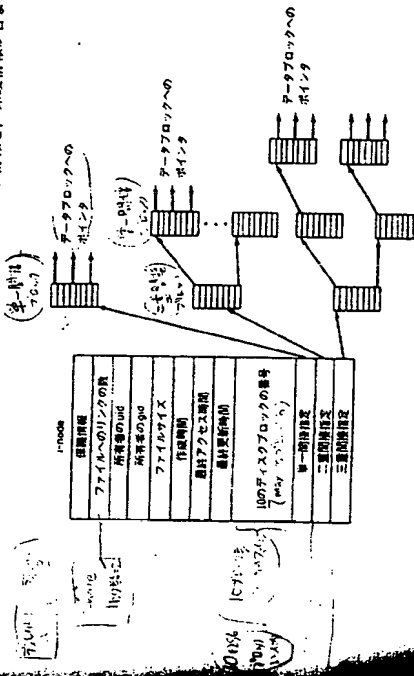


図 5.7 インデックス

る。これらに関しては後述する。ここでは10個のディスクブロック番号と、3つの間接ブロック番号に焦点をあてる。最長10ブロック分のファイルでは、すべてのディスクアドレスが1ノードに保存されているため、ブロックの抽出は簡単である。

ファイルが10ディレクトリ以上の大きな場合、空きディレクトリを確保し、単一階層ブロック内のポイントがそれぞれを指すよう設定される。単一階層ブロックは、ディレクトリブロックに対するポイントを保存するために用いられる。1Kのディレクトリブロックにおいて、32ビットのディレクトリアドレスを用いると、単一階層ブロックでは256個のディレクトリアドレスを保存することができる。この手法により、ファイルは最大256のブロックを持つことができる(1ノード内に10個、単一階層ブロックに256個)。

ファイルの大きさが266ブロック以上になると、二重間接ブロック内のポインタが、最大256のポインタを持つディスクブロックを指すようになる。ただしこれらのポインタはプリアンプブロックを指していない。256の単一間接ブロックを指すのである。二重間接ブロックは、ファイルの大きさが266 + (256 * 65802) = 65802768となるまで使用できる。64M以上の大きさとなるファイルに関しては、三重間接ポインタを用いて、256の二重間接ブロックを指すポインタを含んだブロックを指す。

16G バイトより長いファイルを取扱うことはできない。もちろんディスクブロックの大きさを 2K にすれば、ポインタを持つ各ブロックは、256 ではなく 512 個のポインタを持つことになり、最大ファイルサイズは 128G バイトとなる。128G バイトディスク用の FAT サイズは、予想するものとはばかばかしいほどだ。UNIX スキーマの利点は、間接ブロックが必要の場合にのみ使用されるといふところである。10K 以下のファイルの場合、間接ブロックは必要ない。最も長いファイルでさえも、最高 3 回のディスク参照によってファイル内の任意のバイトのディスクアドレスを得ることができる (i ノードを得るためのディスク参照は除く、i ノードはファイルがオープンされた時点でメモリ内に読み込まれ、ファイルがクローズされるまでメモリに保存される)。

MINIXの手法はUNIXのそれと同じであるが、iノードには7つのディスタブブロック番号が保存することはできず、三重間接ブロックも提供されていない点が異なっている。2バイトのディスタブアドレスと、IKブロックを使用した場合、最大362Mのファイルを取ることができ、大半のパーソナルコンピュータに対しては、十分な大きさである。

5.2.3 デイレクトリ構造

ファイルの読み出しを行うためには、それをオープンしなければならぬ。ファイルをオープンすると、オペレーティング・システムはユーザーが提供したパス名を用い、ディレクトリ構造を探し出し、ファイルの読み書きに用いる。パス名をノードに対応付ける時(あるいはそれに相当するような操作)、ディレクトリ・システムがどの段を構造となっているかを考えねばならない。ディレクトリ構造は簡単にものから、かなり高度なものまで様々である。

銅合金ダイレクトリレーの一例として、CP/Mのダイレクトリレーシステム(Golden Pechura, 2)

5.2 ファイルシステムの設計

1996)を見てもよいでしょう。図 58 は CP/M のディレクトリを明示したものである。このシステムでは 1 つのディレクトリが存在しないため、すべてのファイルはこのディレクトリを検索するだけで検出できる。エントリを発見すると、ディスクブロック番号も共に検出できることになる。ディスクブロック番号が、ディレクトリ・エントリ内に保存されているからである。ファイルが 1 つのエントリに格納されない数のブロックを割り当てられる時、そのファイルは追加ディレクトリ・エントリを割り当てられることになる。

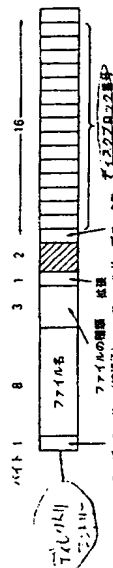


図5.8 ディスクブロック番号を持った各ファイルのディレクトリ

図 5.8 の各フィールドは以下の意味を持っている。ユーザーコードは、どのユーザーがファイルの所有者であるかを記録する部分である。執筆時、現在ログインしているユーザーのエントリだけが対象となる。次の 2 フィールドには、ファイル名と拡張子が格納されている。拡張フィールドは、ファイルの大きさが 16 ブロック以上になった場合に複数のディレクトリ・エントリを割り当てなければならないために必要である。このフィールドには、ディレクトリ・エントリの順番に関する情報格納される。ブロック数は、最大 16 のディスタブロック・エントリのうち、そのいくつかが使用されているかを示す。残り 16 のフィールドにはブロック番号が保存されている。最後のプロックが満杯になっているとは限らないため、システムは最後のバイトまでの正確なサイズを決定することはできない(バイトではなく、ブロック単位でファイルサイズを記録しているため)。

次に階層化されたディレクトリ・ツリーを持つシステムの例を考えてみよう。図 5.9 は MS-DOS のディレクトリ・エントリを示したものである。長さは 32 バイトであり、ファイル名と先頭ブロック番号などが格納されている。先頭ブロック番号は FAT に対するインデックスとして用いられ、それにより 2 番目のブロック番号を導出することになる。2 番目以降のブロック抽出も同様に、それにより 2 番目のブロック番号を導出することになる。

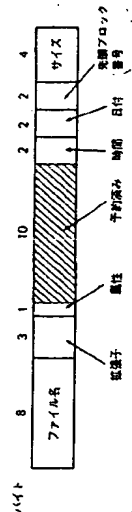


図5.9 MS-DOS/ディレクトリ・エントリ、

5 5 ファイルシステム

様である。この様に、特定のファイルのすべてのブロックを抽出することができる。固定サイズを持つ(360K)のフロッピーディスクでは112エンタリルルートディレクトリを除き、MS-DOSのディレクトリはファイルとして扱われ、任意の数のエンタリルを持たせることができる。

UNIXやMINIXにおけるディレクトリ構造は、図5.10の様に非常に簡単である。各エンタリルはファイル名とそのiノード番号だけを持っている。型、サイズ、時刻、所有者、およびディスクブロックに関する情報はすべてiノード内に格納されている(図5.7参照)。UNIX内のすべてのディレクトリはファイルであり、16バイトのエンタリルを任意の数のだけ持たせることができる。

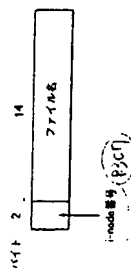


図5.10 UNIXディレクトリ・エンタリ

ファイルをオープンすると、ファイルシステムは与えられたファイル名を用い、ディスクブロックを抽出する。パス名/usr/ast/mboxの抽出がどの様に実行されるかを見てみることにしよう。これはUNIXにおける例であるが、階層化ディレクトリシステムでは同様のアルゴリズムが用いられる。最初にファイルシステムがルートディレクトリを検索する。UNIXではルートのiノードはディスクの決まった場所に置かれている。

次にパスの最初の部分であるusrを抽出するため、ルートディレクトリを検索し、ファイル/usrに対するiノードを得る。このiノードから、システムはusrのディレクトリを検出し、パス名の次の部分であるastを見つける。astのエンタリルを抽出すると、ディレクトリ/usr/astのiノードが得られることになる。このiノードからディレクトリ内容を取り出し、mboxを検索する。このファイルのiノードはメモリに読み込まれ、ファイルがクロージングされるまでメモリ内に置かれる。この検索の過程を図5.11に示している。

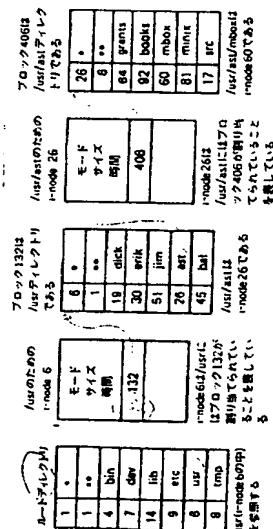


図5.11 /usr/ast/mboxの検索過程

5.2 ファイルシステムの設計

相対パス名の検索も、開始位置がルートではなく、ワーキングディレクトリであるという点を除けば、絶対パス名と同じである。各ディレクトリはディレクトリ作成時に割り当てられた「.」に対するエンタリルを持っている。エンタリル「..」はディレクトリに割り当てられているiノード番号を保持しており、エンタリル「..」は親ディレクトリに割り当てられているiノード番号を保持している。したがって、../disk/progの検索手続きにおいては、単にワーキングディレクトリの「..」を検索し、親ディレクトリのiノード番号を検出し、そのディレクトリ内でdiskを探すことになる。これらのパス名を扱うために、特殊な機構を用いる必要はない。ディレクトリシステムにとって「.」と「..」は、単なる通常のASCII文字列にすぎない。

5.2.4 共有ファイル

何人かのユーザーが同じプロジェクトに携わっている時、ファイルの共有が必要になる。例えばそれぞれのユーザーが所有しているそれぞれのディレクトリから、共有するファイルを同時に読めると便利である。図5.12には、図5.3(c)のファイルシステムを再表示している。ただし今回は、ユーザーCのファイルがユーザーBのディレクトリにも現れている。ユーザーBのディレクトリと共有ファイルの関係をリンク(link)と呼んでいる。ファイルシステム自身は、有向無周期グラフ(DAG: Directed Acyclic Graph)となり、もはやツリーではなくなる。

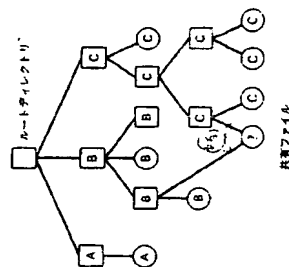


図5.12 共有ファイルを持つファイルシステム

共有ファイルは便利ではあるものの、いくつかの問題がある。まず最初にディレクトリがCPI/Mの様に直接ディスタアドレスを持っている場合である。ファイルがリンクされた場合、ユーザーBのディレクトリにもディスタアドレスを持つことになる。もしユーザーBがユーザーCがこのファイルの後ろにブロックを追加したとすると、新しいブロックはどちらかのユーザーのディレクトリのみには保存されることになる。この変更は別のユーザーからは見ることができなくなり、つまりファイルの共有という目的は果たせないことになる。

5.2 ファイルシステム

この問題は、2つの方法によって解決できる。最初の方法は、ディスクブロックをディレクトリではなく、ファイル自身に対応した小型のデータ構造に持たせることである。これにより、ディレクトリはこのデータ構造だけを指すようになる。この方法はUNIXでも用いられている(その場合、小型のデータ構造に相当するものは「ノード」である)。

もう1つの方法では、システムに新しく「LINK」型のファイルを作成させることにより、ユーザーBがユーザーCのファイルの1つにリンクし、そのファイルをユーザーBのディレクトリに加える。新しいファイルにはリンクされたファイルのパス名だけが含まれている。Bがリンクされたファイルを読み取ると、オペレーティング・システムは読み取られているファイルがLINK型であると判断し、連結されたファイル名を検索し、そのファイルを読み取る。この様な方法はシンボリックリンク(symbolic linking)と呼ばれる。

いずれの方法にも欠点がある。最初の方法において、ユーザーBが共有ファイルにリンクする場合、iノードに登録されたファイルの所有者はCである。リンクによって所有者は変わらない(図5.13参照)。iノードのリンク数は1増やされるため、システムはいくつのディレクトリ・エントリがそのファイルを指しているかが判断できる。

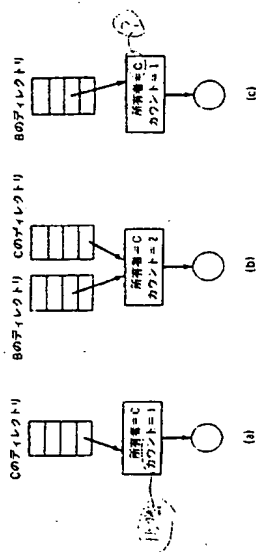


図 5.13 DAG の欠点
(a) リンク前の状態 (b) リンク作成後 (c) 最初の所有者がファイルにリンクした後

その後ユーザーCがそのファイルを削除しようとした時点でシステムに問題が発生する。ファイルを削除し、iノードをクリアすると、Bは無効なiノードを指すディレクトリ・エントリを持つことになる。そのiノードが後に別のファイルに割り当てられた場合、Bのリンクは間違ったファイルを指すことになる。システムはiノード内のカウンタの値から、ファイルがまだ使用中であると判断することはできるが、削除のためにそのファイルに対するすべてのエントリを削除することは不可能である。ディレクトリを指すポインタをiノードに保存することもできない。なぜならディレクトリは無効に存在する可能性があるからである。

そこでiノードは保留し、カウンタを1に設定し、Cのディレクトリ・エントリを除去してしまいうという方法がある。図5.13(c)はこれを示したものである。これにより、ユーザーCが所有して

5.2 ファイルシステムの設計

いるファイルのディレクトリ・エントリを持っているのはユーザーBのみとなる。システムがアバウト(CPU使用時間やディスク使用量などに対する課金)、またはユーザーCがそれを削除するまで、ファイルの料金を徴収されることになる。Bがこのファイルを削除した時点で、カウンタは0になり、ファイルは本当に削除される。

シンボリックリンクでは、本当の所有者だけがiノードを指すポインタを持つことになるため、前述した問題は発生しない。このファイルにリンクしたユーザーは、iノードへのポインタではなく、パス名だけを与えられる。本当の所有者がファイルを削除すると、ファイルは破壊されてしまう。シンボリックリンクによるそのファイルへのアクセスは、システムがそのファイルを検出できなければ失敗に終わる。シンボリックリンクを削除してもファイルにはなんの影響もない。シンボリックリンクの問題点は、余分なオーバーヘッドの発生である。そのパスを含んでいるファイルを読み出し、次にそのパスを解析し、iノードにたどり着くまで、1つ1つ追跡していかなくてはならない。これらの作業においては多くの余分なディスクアクセスが必要である。さらに各シンボリックリンクに対してiノードが1つ余分に必要になり、加えてパスを保存するため追加ディスクブロックも必要になる。ただしパス名が短い場合は、最適化のためにシステムはパス名をiノード内部に保存する。シンボリックリンクにはマシンのネットワーク・アドレスと、そのマシン上のパス名を提供するだけで、世界中のどのマシンのファイルにもリンクできるという利点がある。

またシンボリックリンクやその他のリンクにより、別な問題が発生する。リンクが認めらると、ファイルは複数のパスを持つ可能性がある。特定のディレクトリから始まり、そのディレクトリ内のすべてのファイルとサブディレクトリを検出するというプログラムは、リンクされたファイルを複数回検索することになる。例えばディレクトリとそのサブディレクトリのすべてのファイルをテーマにダウンロードすると、ほとんど同じなダウンロードプログラムでない限り、リンクされたファイルがもとのとおりリンクされるのではなく、ディレクトリに2回コピーされる。

5.2.5 ファイルシステムの信頼性

ファイルシステムの破壊は、コンピュータの破壊よりはるかに悲惨な結果を招くことが多い。例えばコンピュータが火災、落雷、あるいはキーボードの上にこぼしたコーヒーによって破壊されると、費用はかかるが交換品を購入すればことは済む。安価なパーソナルコンピュータの場合、販売店に持ち込めば、数時間で交換してくれるだろう(ただし大学の様に注文書の発行までに、いくつもの委員会の承認を得なければならない場合は別である)。

コンピュータのファイルシステムが、ハードウェア、ソフトウェア、フロッピーディスクをかけたネットワークで、取り返しのつかない状態に陥った場合、すべての情報を復元することは困難であり、時間がかかり、しかもいたいていの場合には不可能である。プログラム、書類、

5章 ファイルシステム

客ファイル、根拠記録、データベース、営業計画、その他のデータを永久に失ってしまった者にとって、その痛手は図り知れないほど大きいものである。ファイルシステムにおいて、装置やメディアに対する物理的な保護を行うことはできないが、情報を保護することは可能である。この項では、ファイルシステムの安全保護に関していくつかの説明を行う。

3章でも説明したように、ディスク内に不良ブロックが存在することも珍しくはない。フロッピーディスクは出荷時に完済状態であっても、使用中に不良ブロックが発生するかもしれない。ハードディスクでは、初めから不良ブロックが存在していることが多い。不良ブロックのないディスクを製造するのは困難である。事実、大半のハードディスク製造元は各ドライブに、検査中に発見した不良ブロックリストを付けている。

不良ブロックに対しては2つの解決方法が取られている。1つはハードウェアの面から、そしてもう1つはソフトウェアの面からの解決法である。ハードウェアによる解決法は、ディスク上の1つのセクタを不良ブロックリストに登録し、交換プロセッサに使用する方法である。コントローラが最初に初期化された時、不良ブロックリストを読み出し、交代プロセッサまたはトラップを不良ブロックと置き換え、その内容を不良ブロックリストに登録する。これによって不良ブロックに対するすべての要求では、交代プロセッサが用いられることになる。

ソフトウェアによる解決法では、ユーザーまたはファイルシステムがすべての不良ブロックを含んだファイルを入念に作成しなくてはならない。これによって不良ブロックは空きブロックから除外され、データファイルに現れることはない。不良ブロックからファイルが読み書きされない限り問題は発生しない。ディスクのバックアップを行う時、誤ってこのファイルを読み出さないよう注意しなければならない。

■ バックアップ

不良ブロックを上記の方法で入念に処理したとしても、ファイルを頻繁にバックアップしておくことは大切である。重要なデータがバックアップが破壊されてから自動的に交代トラップに切り換えるという方法は、競争が激化してから、為小量の損失をかけるようなものである。

フロッピーディスク上のファイルシステムは、そのフロッピーディスク全体をブランクのディスクにコピーするだけでバックアップできる。小型のハードディスクのファイルシステムはディスク全体を、標準の9トラック磁気テープ(1巻あたり50M保存できる)か、ストリーマ・テープ(各種のサイズが提供されている)のいずれかにバックアップすることにより、バックアップできる。

大型のハードディスク(例えば500M程度のもの)になると、ドライブ全体をテープにバックアップするのは時間もかかり、効率が悪い。実行は高速でも容量の半分を無駄にしなければならないが、各コンピュータに2つのドライブを接続するという方法がある。ドライブそれぞれ、データ保存用とバックアップ用に2分して使用するのである。毎晩ドライブ0のデータ部分がドライブ1のバックアップ部分にコピーされ、同時にその逆も行われるといった具合である。図5.14はこれを示したものである。この方法により、たとえば1つのドライブが完全に破壊されても、

5.2 ファイルシステムの設計

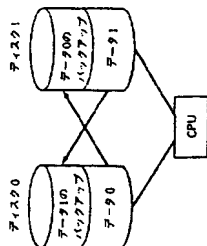


図 5.14 各ドライブを別のドライブにバックアップする
容量の半分の消費されることになる。

情報は保存されることになる。

毎日ファイルシステム全体をダンプする代わりに、増分ダンプ(incremental dumps)という技法が用いられている。増分ダンプの最も簡単な形態は、完全なダンプを週単位や月単位で周期的に行い、この周期的なダンプ以降修正されたファイルのみ日々のダンプを行うというものである。またこれより優れている手法としては、前回ダンプした以降に変更されたファイルのみをダンプするというものがある。

この技法を実現するためには、各ファイルのダンプ時刻のリストをディスク上に保存していないければならない。前回のダンプ以降に変更されている場合は、再度ダンプされ、最終ダンプ時刻が現時刻に変更される。月単位でこれを実行すると、この方法は日々のダンプ用に31本のテープと、月1回行われる完全なダンプを保存できるだけのテープ数が必要となる。その他、これより少ない数のテープを使うという、もっと複雑な手法も用いられている。

■ ファイルシステムの一貫性

もう1つ、信頼性が問題になるのは、ファイルシステムの一貫性である。大半のファイルシステムでは、ブロックを格取り、それを修正し、後でそれを書き込む、修正されたブロックが書き込まれる前にシステムがクラッシュしてしまうと、ファイルシステムは一貫性を失ってしまう。これは書き込まれていないブロックのいくつか、iノードブロック、ディレクトリ・ブロック、あるいは空きリストを含んでいるブロックである場合は特に深刻な問題となる。

一貫性のないファイルシステムを扱うために、大半のコンピュータではファイルシステムの一貫性を検査するユーティリティ・プログラムを持っている。ファイルシステムがアットされた時や、特にクラッシュの後ではこのプログラムを必ず実行する。以下はUNIXやMINIXにおいてこの様なユーティリティがどの様な動きをするかを説明したものであるが、他の多くのシステムでも同じような動きが見られる。ファイルシステム検査プログラムは、各ファイルシステム(ディスク)を個別に検査する。

ブロックとファイルに対して、2種類の一貫性検査が行われる。ブロックの一貫性を検査するため

5.2 ファイルシステム

に、プログラムはブロック2つのカウンタを持つテーブルを作成する。カウンタは当初どちらも0に設定されている。最初のカウンタはファイル内に何度ブロックが現れたかを記録するものである。2番目のカウンタは空きリストに何回現れたかを記録するものである(または空きブロックのビットマップ)。

次にプログラムはすべてのiノードを読み取る。iノードを検査すれば、該当するファイルで使用されているブロック番号のリストを作成できる。各ブロックが読み取られた後、最初のテーブルの各カウンタが増加される。次にプログラムは、空きリスト、またはビットマップを検査し、使用されていないブロックをすべて検出する。空きリストにブロックが見えなくなるたびに、2番目のテーブルのカウンタが増加される。

ファイルシステムに一貫性のある場合は、図5.15(a)に示すように、各ブロックがどこからか、ファイルに1のカウントを持つ。しかしクラッシュが発生すると、テーブルは図5.15(b)の様になる。ここではブロック2はいずれのテーブルにも見えない。これは紛失ブロック(missing block)として報告される。紛失ブロックをのものは特に害はないが、ディスク空間を浪費する。ディスク容量が減少してしまう。紛失ブロックの対処法は簡単である。ファイルシステム検査プログラムがそれを空きリストに入れるだけである。

図5.15(c)の様な状況も考えられる。ここではブロック番号4が、空きリストで2回発生しているのがわかる(空きリストが実際にリストである場合に限り重複が起こる。ビットマップの場合はありえない)。解決策は簡単である。空きリストを再構築すればよい。

最悪な事態として予想されるのは、同じテーブルブロックが図5.15(d)のブロック5の様に、複数のファイルに割り当てられてしまうことである。ファイルのうち、いずれかが削除されると、ブロック5は空きリスト上に配置され、同一ブロックが同時に使用中と、空き状態の両方の状態になってしまう。両方のファイルが削除された場合、そのブロックは空きリストに2回入れられることになる。

ブロック番号															
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	1	0	1	0	1	1	1	0	1	1	0	1	1	0	0
2	1	1	0	1	0	0	1	1	0	0	1	1	0	0	1
3	0	0	1	0	0	0	0	0	0	1	1	0	0	0	1
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
1	1	0	1	0	1	1	1	0	1	1	0	1	1	0	0
2	1	1	0	1	0	0	1	1	0	0	1	1	0	0	1
3	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
1	1	0	1	0	1	1	1	0	1	1	0	1	1	0	0
2	1	1	0	1	0	0	1	1	0	0	1	1	0	0	1
3	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

図 5.15 ファイルシステムの状況
(a)一貫性が保たれた状態 (b)紛失ブロック
(c)空きリストにおける重複データブロック (d)重複データブロック

5.2 ファイルシステムの設計

ファイルシステム検査プログラムは、ここで空きブロックを割り当て、ブロック5の内容をそこにコピーし、そしてコピーしたものをファイルの頭に挿入する。この様にしてファイルの持つ情報はそのままにしておいてもその内容は変わってしまっているだろうか。ファイルシステム構造の一貫性が保たれる。ユーザーが損傷内容を検査する際のために、エラーが報告される。

各ブロックが正しく利用されているかという検査のほかに、ファイルシステム検査プログラムはディレクトリの検査も行う。ここでも2つのカウンタのテーブルを用いるが、ブロックごとではなく、ファイルごとには異なる点がある。ルートディレクトリから始まり、再帰的にツリーを下り、ファイルシステム内の各ディレクトリを検査していく。各ディレクトリ内のファイルごとに、そのファイルのiノードカウンタを増加していく(ディレクトリ・エントリの構造に関しては、図5.10を参照)。

これらの作業が完了すると、iノードを指しているディレクトリ・エントリの数を示した、iノード番号をインデックスとするリストができ上がる。そしてこれらの番号とiノード自身に保存されたリンク数とを比較する。一貫性のあるファイルシステムでは、この2つが一致する。しかしiノード内のリンク数の方が大きい、または小さい場合には、何らかの誤りがある。リンク数がディレクトリ・エントリの数以上の場合は、ディレクトリ内のすべてのファイルが削除されても、カウンタが0以外となり、iノードが削除されないことになる。この誤りはそれほど深刻ではないが、どのディレクトリにも存在しないファイルを持つことになり、ディスクの空間を浪費することになる。iノード内のリンク数を正しい値に設定することにより、これが解決できる。

リンク数がディレクトリ・エントリの数以下である場合は、非常に危険である。2つのディレクトリ・エントリがファイルに連結されたが、iノードでは1つしか存在しないことになっていれば、どちらかのディレクトリ・エントリを削除することによって、iノード内のリンク数は0になってしまふ。i-node内のリンク数が0になった時点で、ファイルシステムはそれを未使用としておき、そのブロックをすべて解放する。この結果、あるディレクトリは未使用のiノードを指してしまうことになり、解放されたブロックは他のファイルに割り当てられてしまう。これも、iノード内のリンク数をディレクトリ・エントリの実際の値に修正すれば、解決できる。

ブロックとディレクトリの検査という2つの操作は、システムの効率を上げるためにはほぼ統合されて行われる(すなわちiノード上では1回のパスだけでよい)。これ以外にも検査を行うことができる。例えばディレクトリは、iノード番号とASCII名を含む特定の形式を持っている。iノード番号がディスク上のiノードの数よりも大きければ、ディレクトリが破壊されていることがわかる。

さらに各iノード内の保護情報型が、0007の様に所有者にもユーザーグループにも全くアクセスを認めず、第3者にはファイルの読み書きや実行を認めるといった、奇妙なものもある。少なくとも、第3者に対して所有者よりも多くの権利を与えているようなファイルを検出させることには意義があるだろう。例えば1,000以上のエントリを持つディレクトリも疑ったほうがよいだ

5.2 ファイルシステム

ろう、ユニザードディレクトリに存在するが、所有者がスーパーユーザーで、SETUIDビットがオンになっているファイルは、その機密保持状態に問題のある場合が多い。少し工夫すれば、奇妙な属性を持つファイルの長いリストにして報告することができる。

以上、ユーザーをクラッシュから保護する際の問題点について説明してきた。ファイルシステムによっては、ユーザーの操作などを考慮したユーザーに対する保護が必要な場合もある。ユーザーが以下の入力を行い、`rm`で終了するすべてのファイル(コンパイルの作成したオブジェクトファイル)を削除しようとした。

```
rm *.o
```

ところが欲って以下の様に人知れずしてしまったとしよう(アスタリスクの後方にスペースがある点に注意せよ)。

```
rm * .o
```

これにより、`rm`はカレントディレクトリ中のすべてのファイルを削除し、`.o`がないというエラーメッセージを表示することになる。MS-DOS、およびあるシステムでは、`rm`を削除すると、単にディレクトリがディレクトリ上に設定されるか、ファイルの削除を示す印がディレクトリに付けられるかの、いずれかである。実際にそれが必要になるまで、ディスクアプロックは空のリストには限られない。したがってユーザーが即座にエラーを発見すると、特殊なユーザーイニテリャープログラムを実行して削除したファイルを取り戻す(復元する)ことができる。

5.2.6 ファイルシステムの性能

ディスクへのアクセスは、メモリへのアクセスよりかなり低速である。メモリの読み出しは通常、多くても数百ナノ秒も要さない。ところがディスクアプロックの読み出しには、10ミリ秒も必要である。メモリの1万倍も、速度が低下する。この様なアクセス時間の違いにより、多くのファイルシステムはディスクアクセス数をできるだけ抑えるように設計されている。

ディスクアクセスを削減するために多用されている技法には、ブロックキャッシュ(block cache)またはバッファキャッシュ(buffer cache)と呼ばれる技法がある(キャッシュはフランス語の *cache* には由来する言葉で、隠蔽の意を持つ)。ここでのキャッシュとは、論理的にはディスクに属しているが、性能上メモリに保存されているアプロックの集合を意味する。

各種のアルゴリズムを用いてキャッシュを管理することができるが、最も一般的なのは、すべての読み出し要求を調べ、必要なアプロックがキャッシュに含まれているかどうかを確認する方法である。含まれている場合には、ディスクにアクセスせずに、読取り要求を満たすことができる。そのアプロックがキャッシュ内に存在しなければ、最初にキャッシュ内にそれを読み取り、そのアプロックを必要とする部分にコピーする。以後、同一アプロックに対する要求があった場合には、キャッシュを用いばよい。

5.2 ファイルシステムの設計

空きがなくなっているキャッシュにアプロックを読み取る場合には、いずれかのアプロックをキャッシュから削除しなければならない。もしキャッシュに置かれてから何らかの変更が加えられている時は、ディスクに書き込まなければならない。これはページングと非常に似ており、4倍で説明した通常のページング・アルゴリズム、すなわちFIFO、セカンダリチャンス、LRUなどがすべて適用される。ページングとキャッシュの主な違いは、キャッシュの参照回数が比較的少なく、すべてのアプロックを正確なLRU順に、連結されたリストとともに保存していくことがそれほど困難ではないという点にある。

残念ながら、ここに落とし穴がある。正確なLRUが可能な状況にあって、LRUを使用することとが望ましくないことがわかった。問題は先の項で触れたクラッシュとファイルシステムの一貫性に関係がある。iノードアプロックの様に重要なアプロックがキャッシュに読み込まれ、修正されたにも関わらず、ディスクに再度書き込まれなかったとすれば、クラッシュによってファイルシステムの一貫性は失われてしまう。iノードアプロックをLRUチェーンの最後に配置したら、先頭にたどり着いてディスクに再度書き込まれるまで、かなりの時間を要することになってしまう。

さらに二重間接アプロックに該当するアプロックでは、短い間隔で2回参照されることは滅多にない。これらを考慮してLRUスキーマを修正してみる。特に、以下の点を配慮する。

- そのアプロックはすぐに再度必要とされるか

- ファイルシステムのの一貫性を保つために、そのアプロックは必要か

これらの質問に対し、アプロックはiノードアプロック、間接アプロック、全体が使用されているアプロック、部分的に使用されているアプロックに分類される。すぐに再度必要とされないアプロックはLRUの終わりではなく、先頭に配置され、そのバッファが即座に再活用されるようにする。書き込まれようとしている部分使用のアプロックなど、すぐに必要となるアプロックは、リストの最後に置かれ、長時間保存されるようにする。

2番目の質問は最初の質問とは関連していない。アプロックがファイルシステムの一貫性を保つために必要なもので(基本的にデータアプロック以外ならすべてこれに当たる)、更新されたものであれば、LRUリストのどこに置かれていても、即座にディスクへ書き込まなければならない。重要なアプロックを即座に書き込むことにより、クラッシュによる被害を最小限にとどめることができる。

この様に、ファイルシステムの一貫性を保つための方法を用いても、実際に書き込むまでデータアプロックを、キャッシュに長時間とどめておくことは望ましくない。パーソナルコンピュータを使って本を書いている人を例にとってみよう。彼が何度もエディタに編集用のファイルをディスクに書き込むよう指示したとしても、まだディスクには何も保存されておらず、キャッシュにすべて残っている可能性がある。システムがここでクラッシュしてしまうと、ファイルシステムの構造は破壊されないが、その日1日の作業内容が失われてしまう。

この様なことが頻繁に発生すると、ユーザーは作業意欲を失ってしまうだろう。この様な問題を

を解決するため、システムは2つの技法を提供している。UNIX ではシステムコール SYNC が用意されており、SYNC では更新されたすべてのブロックを即座に強制的にディスクに書き込む。システムを起動すると、通常 update と呼ばれるプログラムがバックグラウンドで立ち上がり、システムコール SYNC を発行し、次の呼出まで 30 秒間休止するという無限ループを実行する。結果として 30 秒以上間に行われた作業が、グラッシュによって失われることがなくなる。

MS-DOS では、ブロックが更新されると、それをすぐにディスクに書き込んでいく。修正されたすべてのブロックが即座に書き込まれる方式のキャッシュは、ライトスルーキャッシュ(write through cache)と呼ばれている。これは一般のキャッシュよりかなり多くのディスク入出力操作を必要とする。2つの方法の相違は、プログラムが 1K ブロックが満杯となるまで、1文字ずつ書き込んだ後に明らかに明らかなる。UNIX はキャッシュ内にすべての文字を貯めておき、そのブロックを 30 秒に 1 回ずつ、またはブロックがキャッシュから削除されるたびに書き込む。MS-DOS は書き込まれた文字すべてに対してディスクアクセスを行う。もちろん大半のプログラムは内部バッファリングを行い、1文字ずつではなく、行やそれより大きな単位で WRITE システムコールにより書き込みを行う。

この様なキャッシング技法の違いにより、UNIX システムから SYNC を行わずに(フロッピー)ディスクを取り出すと、たいいていの場合データの損失を招くだけでなく、しばしばファイルシステムの破壊にまで至ってしまう。MS-DOS ではこの様な問題は起こらない。この様な異なる方法が用いられている背景には、UNIX は、すべてのディスクが取りはずしできないハードディスクであるという環境で開発され、MS-DOS は、フロッピーディスク環境から始まったという背景がある。ハードディスクがマイクログコンピュータでさえも一般的になるにつれ、より効率的な UNIX の技法が多用されるようになるだろう。

キャッシングだけがファイルシステムの性能を向上させる方法ではない。もう1つ重要な技法として、連続アクセスを行うブロックをできるだけ同じシリンダ内の近接する場所に置くことにより、ディスクアームの移動量を削減する方法がある。出力ファイルが書き込まれると、ファイルシステムはブロックを1つずつ必要に応じて割り当て、空きブロックがビットマップで管理されているれば、直前のブロックに最も近い空きブロックを抽出することは非常に簡単である。空きリストを用いており、その一部がディスクに存在する場合は、近接したブロックの割当てはかなりの困難になる。

しかし空きリストを使っても、ある程度のブロックラスタリングを行うことは可能である。その際はディスクの保存内容をブロックではなく、連続するブロックグループで記録していくのがこつである。1トラックに 512 バイトのセクタが 64 個含まれているとすれば、システムは 1K ブロック(2 セクタ)を使用するが、ディスクの保存容量の割当ては 2 ブロック単位(4 セクタ)で行う。これは 2K のディスクブロックを使用するのとは異なる。キャッシュも 1K のブロックを用いて、ディスク転送も引き続き 1K 単位で行われるが、ファイルをアイドル状態のシステムで連続的に読み出す場合は、シーク回数は 2 分の 1 に減少し、性能をかなり向上させることができる。

同じ技法に若干変化を加え、回転配置を考慮したものがある。ブロックを割り当てる時、システムはファイル内の連続ブロックを同じシリンダに配置しようとするが、最大スループットを得るために化ラングリーになる。したがってディスクの回転時間が最大 16.67 ミリ秒であり、ユーザープロセスがディスクブロックを要求してから、それを入力するまで 4 ミリ秒かかる場合は、各ブロックは少なくとも前のブロックから 1/4 の距離をおいて配置されなくてはならない。

さらに、i ノードもしくは、それに相当するものを用いているシステムにおいて、性能向上を妨げるのは、短いファイルを読み出す場合でも i ノードとブロックに対して合計 2 回のディスクアクセスが必要になるという点である。通常の i ノード配置を示したのが図 5.16(a) である。ここではすべての i ノードがディスクの先頭近くに配置されており、i ノードとブロックの平均距離がシリンダの約半分であるため、シーク距離も長くなる。

ここで i ノードをディスクの先頭ではなく、中央に置き、i ノードと最初のブロックとの距離を 2 分の 1 に削減することによって、容易に性能を向上させることができる。図 5.16(b) に示されたもう 1 つの方法は、ディスクをシリンダグループに分割し、それぞれのグループに i ノード、ブロック、空きリストを持たせるという技法である(McKusick 他, 1984)。新しいファイルを作成する時、任意の i ノードを選択することはできるが、その際、i ノードと同じシリンダグループ内にブロックを格納する、使用可能なブロックがなければ、近接のシリンダグループ・ブロックが使用される。

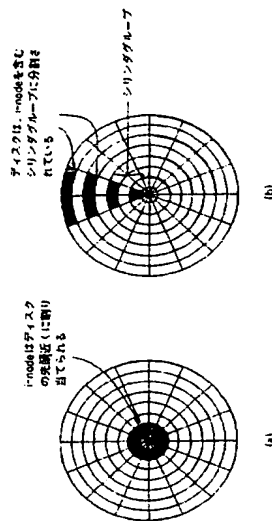


図 5.16 i ノードの位置
(a) ディスクの先頭に置かれた i ノード。
(b) シリンダグループに分割されたディスク。それぞれが独自のブロックと i ノードを持っている。

5.3 ファイルサーバー

分散処理システムには、他のマシンにファイル機能を提供するマシンが存在していることが多い。これらのマシンはファイルサーバー (file servers) と呼ばれている (Birrell Needham, 1980; Fridrich Older, 1981; Svoboda, 1984; Swinehart 他, 1979)。分散処理システムのコストを抑えるための、一般的な方法として、ユーザーにディスクレス・ワークステーションを与え、共通のファイルサーバーに READ および WRITE 要求を送ることにより、ネットワーク経由でそのファイルにアクセスするという方法がある。

原則的に、この様なファイルサーバーは、UNIX、MS-DOS あるいは他の一般的なファイルシステムと全く同じインターフェイスを持つことができる。実際には研究者がいろいろな新しいアイデアをばねし構築した。結果として多くのファイルサーバーがこれまでのファイルシステムには見られなかった機能を持つようになった。以下の項ではそのうちいくつかを見ていくことにする。

5.3.1 インターフェイス・レベル

ファイルサーバーは以下の3つのレベルのいずれにおいてもユーザー(クライアント)にインターフェイスを提供することができる。最も簡単なものは遠隔ディスク(remote disk)である。このモデルでは、各ユーザーに対してファイルサーバーのディスク上に専有部分を仮想ディスク(virtual disk)として割り当てる。ユーザーはローカルディスクと同様に仮想ディスクを扱うことができる。ファイルサーバーはコマンド READ BLOCK と WRITE BLOCK を、ローカルディスクと同じように提供する。

次は、ファイルサーバーはファイルだけを待ち、ディレクトリは持たないという方法である。ファイルの作成、削除、読み出し、書き込み、およびシークを行うコマンドが提供されている。ユーザーがファイルを作成すると、ファイルサーバーは以降の操作に用いる識別子を返して、識別子には例えばランダムな長い番号など、不正ユーザーが識別子に似るものを用いる。ASCII 名をファイルサーバーの識別子に関連付けるためのディレクトリを返せるか否かは、ユーザー自身が決める。この様なディレクトリを設ける場合は、識別子には UNIX のディレクトリに保存されている i ノード番号に類似したものを用いることもできる。

このスキームにおける1つの問題は、ユーザーがファイルサーバー上にファイルを作成し、ディレクトリに識別子を書き出す前にシステムがクラッシュした場合、ファイルが「紛失」してしまう点にある。永入に存在するが、識別子がわからなくなためアクセスができなくなる。この様な状態から脱出する唯一の方法は、ユーザーが自身のファイルの全リストを要求できるようなコマンドをファイルサーバーに設けることである。

3番目のインターフェイス・レベルは、ファイルサーバーに UNIX の様な完全なファイルシステム

システムを持たせつつ、さらにより高度な機能を実現するという方法である。この方法を用いると、コマンドはファイルの操作だけでなく、ディレクトリの作成、削除、ワーキングディレクトリの変更、既存ファイルへのリンク作成と解除なども行えるようになる。リモート・ファイルサーバーが完全なディレクトリ機能を提供する時、ユーザーのワークステーションがリモート・ファイルシステムをワークステーションのファイルツリーにマウントできることもある。その様な場合、マウント後はリモートファイルへのアクセスが、ローカルマシン上のルートからの絶対パスを指定するだけで行えるようになる。したがってそれがリモート・ファイルシステムであるということ意識する必要もなくなる。

5.3.2 アトミック更新

高度な性能が要求されたファイルシステムは、結果としてある意味ではかつてのものほど信頼性がないことが判明した。その昔、まだディスクが突如として壊れる時代には、企業がたいては在庫管理方法を考えてみよう。通常マスターテープに全製品のリストと、その在庫数を記載していた。たいていはオリジナルのテープがほこりや、湿気、テープドライブによる破損によって万が一読取不能になった場合を考慮して、バックアップ用のコピーがとられ、同じように保管されていた。

毎日一度マスターテープがドライブ1に装着され、その日の売上を記録したテープがドライブ2に、そしてクラッシュテープがドライブ3に装着された。そして更新プログラムを実行し、マスターテープと売上テープを読み出し、新しいマスターテープを作成した(そしてバックアップ用のコピーを取られた)。このプログラムが途中でクラッシュしてしまっても、3つのテープを巻き戻して、最初からやり直せばよかった。

このシステムには更新が問題なく完了しても、失敗に終わっても、オリジナルのテープには元の状態も及ばないという利点がある。さらにマスターテープが破損しても、いつでもバックアップ・コピーを用いることができる。

ディスクが発明されると、当然のことながらマスターテープは、ディスク上のファイルとなった。マスターテープの更新は、このファイルを読み出して更新する作業に変わった。唯一の問題は、更新中にクラッシュしてしまった場合、システムを復旧し、もとの状態に戻し、更新プログラムを再実行することが不可能な点であった。すでにどれだけのレコードが更新されているか判断できないからである。

3つのテープを用いたシステムには見られなかったが、ディスクを用いた更新においては見られなかった性質は、アトミック更新(atomic update)または故障復元性(failure atomicity)と呼ばれている。すなわちレコードに対する更新は安全に成功するか、しない場合にはシステムはもとの状態のいずれかである。アトミック更新が失敗に終われば、更新プログラムが再度実行されるだけである。回避しなくてはならないのは、データの全部ではなく、一部だけを更新するような

更新を行うことにより、ファイルの一部を部分的に更新しただけの、不明瞭な状態にしてしまうことがある。

テープシステムの性質でもう1つ注目すべき点は耐故障性(fault tolerance)と呼ばれるものがある。マスタータープに接続できない場所があっても、バックアップをとってあるため心配する必要はなかった。理論的にはファイルシステムは、すべてのファイルのコピーを2つ保存することができると、実際に2つコピーをとるシステムはほとんどない。

耐故障性を提供しているファイルサーバーでは、通常図5.17の様にアトミック更新によって、2つの物理ドライブから論理ドライブを実現することになる。論理ブロックnに情報が書き込まれた後、まずサーバーはそれをドライブ1上の物理ブロックnに書き込み、そしてそれをドライブ2上の物理ブロックnに書き込み、それを検証する。

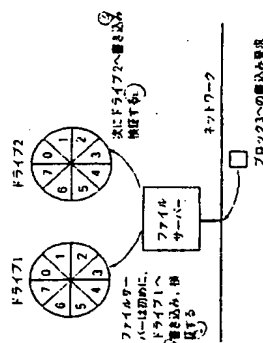


図 5.17 ファイルサーバーの耐故障性の実現
同じ情報を持つ2つのディスクドライブを使用することによって、安定した保存を行うことができる。

この技法は安定記憶機構(stable storage)と呼ばれており(Lampson Sturgis, 1979)、興味深い性質を持っている。まず不良ブロックが偶然どちらかのブロックにできたとしても、周期的に、例えば毎晩ファイルサーバーは不良ブロック、すなわちチェックサムエラーを持つブロックを検出するために、両方のディスクを読み取る。不良ブロックを検出した場合、他のドライブのコピーを、不良ブロックに上書きする。それによってエラーを回避することができる(ブロックが物理的に破損している時、サーバーはそのためにあらかじめおいた交代トラックを用いることができる)。2つのブロックが同時に不良になってしまいう可能性が非常に低いと仮定すると、頻りにディスクエラーが発生したとしても、安定記憶機構によるデータ損失は考えられない。

それではサーバークラッシュの影響について考えてみよう。どちらかのドライブに書き込みを行っている間にクラッシュが発生すると、書き込まれているブロックはチェックサムエラーを引き起こす。このエラーは必ず検出されるため、正常なブロックを使って不良ブロックを置き換え

るとよい。サーバーがドライブ1に書き込みしている間にクラッシュした場合、システムはもとの状態に復元される。ドライブ2に書き込みしている最中にクラッシュすると、システムは更新された状態に復元される。安定記憶機構に対する書き込みが成功しようと、失敗しようと、あいまいな状態になってしまいうことはない。

アトミック更新にやや関連しているのがマルチバージョンファイル(multiversion files)である。ファイルサーバーがマルチバージョンファイルをサポートしている場合、いったん作成された後には、二度と再びファイルに対する更新は行われない。その代わりにファイルの一時的なコピーを作成し、それをアトミック更新によって変更し、そのファイルを最新のバージョンとする。したがってファイルは時系列の不変のバージョンから構成されている。ファイルの読取り要求では、特に明示されていない限り、必ず最新のバージョンを読み出す。

5.3.3 同時実行制御

典型的なファイルシステムでは(例えばUNIX)、2人のユーザーが同時に1つのファイルを更新すると、コマンドREADとWRITEは到着順に実行される。しかし例えば銀行で、2人の顧客がそれぞれ同じ口座に同時に振込みを行おうとしたらどうなるだろう。その口座には当初500ドルとなっており、顧客はそれぞれ200ドルと300ドルを振り込もうとしている。以下の様な取引内容となるだろう。

- ①顧客1のプログラムは残高を読み出し、500ドルであることを知る。
- ②顧客2のプログラムは残高を読み出し、500ドルであることを知る。
- ③顧客1のプログラムは残高を500+200=700に更新する。
- ④顧客2のプログラムは残高を500+300=800に更新する。

最終的な残高は800ドルになる。顧客1が顧客2より遅く処理されれば、最終的な残高は700ドルになっていたであろう。どちらにせよ2つの更新が同時に行われているため、最終的な結果は誤ったものとなる。ここで必要なのは、最初の顧客がプログラムを実行し、次にもう1人の顧客が実行するといった方法で、上記の様な状態が発生しないようにすることである。

同時に更新が行われても、任意の順序で更新を行った場合と同じ結果となるようにすることを直列化(serializability)と呼んでいる。また直列化のための技法を同時実行制御アルゴリズム(currency control algorithms)と呼んでいる。アトミック更新と同じように、データベースシステムではすでに普及しており、今後はファイルサーバーにも採用されることになるだろう。

多くのファイルサーバーが同時実行制御技法として、なんらかの形でロック(locking)機能を提供している。あるクライアントがファイルをロックすると、他のクライアントがそのファイルを使用、またはロックしようとしてもサーバーが拒否する。ロックを使い、次の様な方法で同時実行制御を行う。更新を実行する前に必要なファイルをすべてロックする。もしいずれかのファイルがすでにロックされている(他のクライアントによって)、ロックしたファイル

5章 ファイルシステム

をすべてロック解除し、何も変更されない状態で実行が失敗に終わる、すべてのファイルのロックが成功した場合には、ファイルの読み書きが行われ、ロックは解除される。

ロックアップにおける重要な問題は、クライアントがファイルに対してロックを要求し、そのあとクラッシュしてしまったらどうなるのかという点である。あるサーバーではロックされたままのファイルが復元するようにするために、ロックが行われると必ずタイムアウトを起動し、ロックが解除される前にタイムアウトが切れると、サーバーはクライアントがクラッシュしたものとみなし、ファイルのロックを解除する。しかしこの様な方法では、もしクライアントが単に低速で作業を行っているのみであり、異常がない場合には不都合が生じる。

同時実行制御は2章で学習した相互排他問題に非常に似ているが、若干異なっている。2章ではプログラムの側から問題を考え、危険領域を設定することによって、2つの危険領域が同時に使用されないようにした。ロックではデータの側からこの問題をとらえ、どのプログラムが実行中であるかに関係なく、各ファイルと直接関係していた。特定のファイルにアクセスする可能性のあるプログラムが多数存在する時、プログラムではなく、ファイルに制御を行ったほうがよいだろう。

実現方法ももちろん異なっている。なぜなら、ファイルシステム内にそれぞれファイルに付するセマフォを持ち(まれにしか発生しない)、ロックに備えらる、無敵が多くなる。その代わりロックされたファイルの一覧がメモリに保存される。また、クラッシュ、時間切れ、などのこのほかの部分も異なっている。

5.3.4 トランザクション

自動ロッキングは、トランザクション(transaction)という形でアトミック更新と組み合わされることもある。トランザクションは、成功すれば最後まで実行され、失敗すればシステムには人の影響を与えないという特徴を持っている。トランザクションを実行するには、クライアントプロセスがメッセージ "BEGIN TRANSACTION" をファイルサーバーに対して送信しなければならない。そして必要な数のファイルを読み取る。作業が完了すると、プロセスはメッセージ "END TRANSACTION" をファイルサーバーに送り、すべての変更内容をコミット済み(committed) すなわち1つのアトミック更新において不変のものとする。コミットが不可能な場合は、トランザクションは失敗に終わり、なんの変化も起こらない。メッセージ "END TRANSACTION" を受け取るまでは、ファイルを使用している他のユーザーが書き込まれた内容を見ることができない。他のユーザーが参照できるのは、修正前のファイルだけである。

複数のプロセスが同時にトランザクションを実行している時、ファイルサーバーはその時点で実行中であるのはそのプロセスのみであるという認識を起こさせる。つまり、ファイルサーバーは自動的にすべてのトランザクションを直列化しなければならない。ファイルサーバーが、直列化とアトミック更新を複数のファイルに対して効果的に行うための方法を見つけたのは、今後の課題である。ある簡単な方式を大まかに示すと、次の様になる。

5.3 ファイルサーバー

プロセスがトランザクションを開始すると、ファイルサーバーは安定記憶機構上にトランザクションレコード(transaction record)を作成し、その状態を記録する。安定記憶機構上に格納されているため、サーバー、ディスクのいずれがクラッシュしても失われることはない。プロセスが初めてファイルを読み取る際、ファイルサーバーはファイルをロックして他のプロセスのアクセスを防止する。ロックを行うことができれば、トランザクションが失敗に終わる、何も変化は起こらない。

プロセスが初めてファイル上にブロックを書き込もうとすると、ファイルはロックされ、ファイルのコピーがとられる。書き込みはオリジナルでなく、このコピーに対して行われる(そしてそのファイルに対する以降の書き込みもすべてコピーに対して行われる)。

メッセージ "END TRANSACTION" が実行されると、サーバーは新しく作成されたファイルの集合と、まだ修正されていない古いファイル集合を置き換える。そしてインテンションリスト(intentions list)を作成し、更新が必要なファイルと、それぞれの新しいファイル位置を表示する。インテンションリストは安定記憶機構のトランザクションレコードに加えられる。サーバーがクラッシュしても、再度立ち上げた時点で復旧できるようにする。そしてトランザクションが完了し、コミット済みと印付けするまで、他の操作は行えないことになる。

次に上書きされるファイル内のブロックをすべてについて、そのレコードをアトミック更新することにより、新しいファイルと置換する。この時点でトランザクションは完了し、すべてのロックが解除される。これらの手順を図5.18にまとめている。

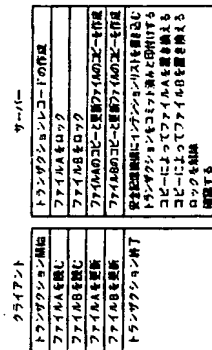


図5.18 トランザクションの手順

- (a) トランザクション中にクライアントが行う手順
(b) サーバー側においてサーバーが行う手順

これ以外により効率の良いアルゴリズムが発表されている。詳細はSvobodova (1984) の論文と、その参考文献を参照されたい。

5.3.5 複製ファイル

ファイルサーバーによっては、複製(replication)と呼ぶ便利な機能を用いているものがある。各ファイルを1つだけ保存するのではなく、サーバーは n 個のファイルを保存する。コピーの1つが壊れて破壊されても、データが失われることは決してない。標準ディレクトリが(ASCII 名, i ノード番号)のペアとしてリストされると、ディレクトリを n 個のファイルの(ASCII 名, n 個のiノード番号)を含むように変更し、複製を行うことができる。これを示したのが図5.19である。複製ファイルの詳細はPope他(1986)による文献などを参照されたい。

ASCII 名		inode	
ファイル1	17	17	40
ファイル2	22	22	91
ファイル3	12	30	29
ファイル4	84	15	66

(a)

図5.19 複製ファイル
(a)通常のディレクトリ (b)複製ディレクトリを持つディレクトリ

ファイルの1つのコピーが修正された時、その複製に対する操作はファイルサーバーが行う。2つの方法が考えられる。最初の方法はディレクトリにリストされている複製に対して、変更したブロックをそれぞれに送り、更新を行う。もう1つの方法はすでに古くなってしまった複製ファイルを削除して、修正済みのファイルの新しいコピーを作成し、これらをディレクトリに配置するという方法である。

異なるファイルサーバー上に複製ファイルが保存され、そのうちいくつかは複製作業中のネットワークのダウンにより、孤立してしまっても、分散されている他のコピーがバックアップ時に用いられることになる。複製ファイルがたまたまディレクトリであり、ネットワークがダウンしている間にすべてのコピーに対して個別に修正が行われた場合、ネットワークが再度立ち上がった時にファイルサーバーは一貫性を失った複数のディレクトリ:コピーを持つことになる。そこで、一貫性を取り戻すためになんらかの策を講じなくてはならない。この問題の詳細とその解決策についてはPope他(1981)、Walker他(1983)、Weinstein他(1985)などの論文を参照されたい。

5.4 セキュリティ

ファイルシステムには、ユーザーにとって非常に重要な情報が含まれていることが多い。したがって、この情報を不正なユーザーから保護することはファイルシステムにとって最も大切なことである。以下の項ではセキュリティと保護に関する各種の問題について見ていくことにしよう。

5.4.1 セキュリティ環境

「セキュリティ」と「保護」という言葉は、同じ意味を持つ言葉として使われることも多いが、次の様な2つの一般的な問題として区別されている。ファイルが不正なユーザーによって読み書きされてはならないということから、技術的、管理的、法的、および政治的な要因も含まれている。もう1つはセキュリティを提供するための特定のオペレーティング・システムの機構である。本書では混乱を避けるために、全体的な問題を指す時にはセキュリティ(security)という言葉を用い、そしてコンピュータ内の情報を守るために用いられる特定のオペレーティング・システムの機構を指す時には保護機構(protection mechanisms)という言葉を用いることにする。ただしこの2つの境界はあまりはっきりと定義されていない。最初にセキュリティについて説明し、本章の後半において保護機構を見ていくことにしよう。

セキュリティには、多くの側面がある。そのうち最も重要なものを2つ挙げるとすれば、データの破壊と侵入であらう。データ破壊の主な原因としては、次の様なものが挙げられる。

- 天災: 火災、洪水、地震、戦争、暴動、ネズミがカード、テープ、フロッピーディスクをかじったなど。
- ハードウェアまたはソフトウェア上のエラー: CPUの誤動作、読取り不能なディスクまたはテープ、通信エラー、プログラムのバグ。
- 人為的エラー: データの誤入力、誤ったテープやディスクのマウント、誤ったプログラムの実行、ディスクまたはテープの紛失。

このうち大部分はオリジナル・データからできるだけ離れたところに、適切なバックアップを保存することによって回避できる。

より興味深いのは、侵入問題の対処方法である。侵入には2つの形態がある。受動的な侵入は、読取り許可を与えられていないファイルの読取りである。能動的な侵入は、より悪意が強いものである。不正な変更を意図的に行うとするとした場合がこれに相当する。システム設計において侵入に対するシステムの保護をより強力にするためには、どの様な侵入に対して保護を行うべきかを理解する必要がある。一般的な侵入方法には次の様なものがある。

5章 ファイルシステム

●技術者以外の悪意、根上にタイムシェアリングシステムの端末を持っているユーザーは多く、近接したユーザーの中には、特に障害がなければ他のユーザーの電子メールやファイルの内容を見ようとする者もあるだろう。例えば大抵のUNIX システムでは特に制限を与えない限り、すべてのファイルを誰でも読み取ることができる。

●内部の人間による悪意、学生、システムプログラマー、オペレータ、あるいはその他の技術者によっては、コンピュータシステムのセキュリティを破ることに意欲を燃している者もある。高度な技術を持ち、かなりの時間を費やして目的を達しようとすることもある。

●一攫千金という人、銀行からお金を盗むために銀行システムに侵入しようとした、銀行のプログラマーがいた。ソフトウェアに修正を加え、四捨五入ではなく、切り捨てを行わせたり、セント以下を切り捨てることによって利益を得たり、何年間も使用されていない口座から預金を吸収したり、脅迫状を送ったりと(「要求に応じなければ、銀行の取引記録をすべて破壊してしまう」などと脅す)。その方法は様々である。

●産業または軍事スパイ、御会企業または外国など資金力のある団体が行うスパイ行為、プログラム、企業秘密、特許、技術、回線設計、営業プランなどを盗もうとする行為である。この様な行為は、しばしば盗聴や、コンピュータからの電磁放射を捕えるためのアンテナを使ったりする。

KCBから軍事機密を守らうとする場合と、学生が奇妙なノックセンサをマシンに不正に挿入するのを防ぐ場合では事の重大さが違う。セキュリティと保護にかけける意気込みは、誰が敵であるかによって明らかに違ってくる。

セキュリティのもう一つの側面は、プライバシー(privacy)である。すなわち個人情報の取扱いを防ぐことが必要である。これは多くの法制的および道義的問題を引き起こすことになる。もし政府が独自のやり方で各個人の身の上書を作成し、社会保障の費用や税金をこまかした人を洗いだそうとしたらどうなるだろう。警察は組織犯罪を防止するために、誰のどの様な情報でも見ることが出来るのだろうか。雇用者や保険会社はその様な権利を持っているだろうか。これらの権利と個人の権利が矛盾する時はどうすればよいのだろうか。これらの問題はすべて非常に重要であるが、本書では目及しないことにする(参考文献に関しては APPENDIX E を参照されたい)。

5.4.2 セキュリティの落とし穴

運搬業界でのタイタニックやセンデンプラダの様に、コンピュータの機密保護の専門家にも思いついたくない事がある。この項では、4つのオペレーティングシステム(UNIX, MULTICS, TENEX, およびOS/360)で発生したいくつかのセキュリティの問題を見ていくことにする。

UNIX のユーザーリティである `lpr` は、ラインプリンタにファイルを印刷するものである。これ

5.4 セキュリティ

には印刷後ファイルを削除するというオプションが提供されている。UNIX の初期バージョンでは、誰もが `lpr` を使ってパスワード・ファイルを印刷し、その後システムにパスワード・ファイルを削除させることができた。

UNIX に侵入するもう一つの方法は、ワーキングディレクトリ内に存在する `core` という名前のファイルをパスワード・ファイルにリンクするものである。そして侵入者が `SETUID` を持ったプログラムを実行した際、コアダンプを強制的に引き起こすことにより、システムはパスワード・ファイルの先頭から `core` ファイルを書き込む。こうしてユーザーはパスワード・ファイルを、自分の選んだいくつかの文字列(例えばコマンド引数など)を持つファイルに置き換えることができる。

さらに UNIX には以下のコマンドに関わる問題もある。

```
mkdir foo
```

`mkdir` は、`root` が所有者であり、`SETUID` を持つプログラムである。システムコール `MKNOD` によってディレクトリ `foo` の `i` ノードを作成後、`foo` の所有者を有効 `uid` (すなわち `root`) から、実 `uid` (ユーザーの `uid`) に変更する。システムが迅速で動作している時、ユーザーはディレクトリ `i` ノードを素早く削除し、`MKNOD` と `CHOWN` の間に `foo` という名前がパスワード・ファイルに對するリンクを作成できることがある。 `mkdir` が `CHOWN` を実行すると、そのユーザーがパスワード・ファイルの所有者になる。必要なコマンドをシェルスクリプトとして用意することで、リンクが成功するまで何度もやり直すことができる。

MULTICS におけるセキュリティの問題は、システム設計者が MULTICS をタイムシェアリングシステムとして考えており、高価なバッチ処理者のために後から思い付きでバッチ機能を追加したという点に起因している。タイムシェアリングのセキュリティは優れていた。しかしその一方ではバッチシステムにおけるセキュリティが存在しなかった。任意のユーザーディレクトリに多数のカードを読み出させるようなバッチジョブを挿入することぐらいであれば、誰かファイルを読み出すためには、エディタのソースコードを入手し、ファイルを読むよう修正した(ただしエディタの機能はそのままで、それを犠牲者のディレクトリ `bin` に置いておく。その犠牲者がエディタを呼び出すと、侵入者の作成したエディタが実行される。エディタは普通に動くが、犠牲者のファイルが読まれてしまう。通常のプログラムを修正し、通常の機能以外に悪意の機能を持たせ、犠牲者に修正したものを使用させるという技法はトロイの木馬 (Trojan horse attack) と呼ばれている。

TENEX オペレーティングシステムは、かつて DEC-10 コンピュータ上で、人気を博していた。今ではそれはほとんど人気はなくなったが、次の様な設計上の誤りにより、コンピュータ・セキュリティの歴史に永久に残ることになってしまった。TENEX もペーシングをサポートしていた。ユーザーは自身のプログラムの動作を監視するために、ページフォルト発生のためにユーザー関

5.4 ファイルシステム

教を呼び出すようシステムに指示することができた。

TENEX ではパスワードを用いてファイルを保護することも可能であった。ファイルにアクセスするために、プログラムは正しいパスワードを提示しなければならない。オペレーティングシステムは1文字ずつパスワードを検査し、誤りを発見した場合に、ただちに停止する。TENEX に投入するためには、図5.20(a)に示すように、最初の文字を1ページの最後に配置し、残りを次のページの最初に書くといった操作を注意深く行う。

次の手順は、例えば非常に多くのページを参照することにより、2ページ目をメモリから追い出し、3ページ目をメモリ内に置かないようにすることである。ここでプログラムは用意されたパスワードを用いて、権限者のファイルをオープンしようとする。実際のパスワードの最初の文字がA以外ならば、システムは最初の文字の検査で停止し、ILLEGAL PASSWORDを返す。しかし実際のパスワードがAから始まる場合、システムは検査を続け、ページフォルトに遭遇した時点で侵入者に通知される。

パスワードがAで始まらない場合は、侵入者は図5.20(b)の様なパスワードに変更し、プロセス全体を繰り返して、Bから開始するかを調べる。最高128回のプロセスを繰り返して、ASCII文字のすべてを調べれば、最初の文字が判明する。

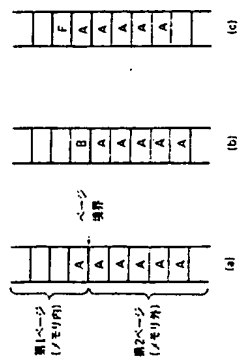


図 5.20 TENEX パスワードの問題点

最初の文字がFだったとしよう。図5.20(c)のメモリアクセスで侵入者はFA、FBといった文字列の検査を行うことができる。この方法を用いることにより、n文字のパスワードを見つけてするために必要な試行回数は、 128^n ではなく、 $128 \times n$ 回となる。

最後の問題はOS/360に関するものである。以下は概略であるが、問題の基本的な部分は解説している。このシステムではデータの読取りを開始し、ドライブがデータをユーザー空間に転送している間、演算を行うことが可能であった。そのトリックとは、まず任意深くデータの読取りを開始し、次にユーザーのデータ構造体(例えば読み取るファイル名とそのパスワード)を必要とするシステムコールを発行することである。

このオペレーティングシステムは、最初にパスワードが指定したファイルに対して正しいもの

5.4 セキュリティ

であることを確認する。次に、実際にアクセスを行うために再度パスワードを取り出すとする(パスワードをシステム内部に保存しておけばよかったのだが、実際には行っていなかった)。残念ながら、システムが2回目にパスワードを取り出す直前に、パスワードはテープドライブからのデータによって上書きされている。そしてシステムはパスワードの検査を行わずに新しいパスワードを読み出してしまふ。この様にタイミングをうまくつかむにはある程度経験が必要だが、それほど難しいことではなかった。しかもコンピュータが得意とする、同一操作の無数の繰り返しが含まれていた。

5.4.3 一般的なセキュリティへの攻撃

前述の様な問題はすでに改善されているが、平明なオペレーティングシステムではまるまるいかに情報を漏れさせているのが実状である。システムのセキュリティ状態を知る一般的な方法としては、タイガーチーム(tiger teams)や、透過チーム(penetration teams)と呼ばれる専門家グループを雇い、侵入が可能かを調べてもらう。Hebbard(他 1980)も同じことを大学院生を使って行った。透過チームは何年かの期間をかけ、システムを弱点と考えられる部分をいくつか発見した。以下に成功する確立の高い攻撃方法をいくつか挙げてみた。システムを設計する際、これらの点には特に注意されたい。

- メモリページ、ディスク空間、テープを要求し、それを読み取る。多くのシステムは割当て前にそれを消去することはない。前の所有者が書き込んだ興味深い情報で満たされているかも知れない。
- 不正なシステムコール、不正な引数を持つシステムコール、または不正ではないが意味のない引数を持つシステムコールの使用。多くのシステムは混乱を生じる。
- ログインを開始し、DEL, RUBOUT, BREAKをログインセッションの途中で入力する。システムによってはパスワード検査プログラムが強制終了させられ、ログインが成功したとみなされる。
- ユーザー空間に保存されている複雑なオペレーティングシステムの構造体を修正してみる。多くのシステムではファイルをオープンするために、プログラムが大きなデータ構造を用意し、それにファイル名やその他の多くの引数を設定し、システムに渡す。ファイルの読み書きを行うと、システムは構造体自身も更新することがある。これらのフィールドを変更するとセキュリティに大きな障害が出てくる。
- 画面に"login:"を表示するプログラムを作成し、ユーザーを混乱させる。ユーザーの多くは画面に向かってログイン名と、パスワードを入力してしまい、プログラムはそれを投入のために記録する。

5 章 ファイルシステム

● マニュアルで「～してはならない」と記されている事項を採し、できるだけ多くの禁止事項を行ってみる。

● システムプログラムを脱得し、特定のログイン名を持つユーザーに対して重要なセキュリティ検査のいくつかをとばすようシステム内容を変更させる。この攻撃方法は落し穴(trapdoor)技法として知られている。

● 上述のいかなる方法も使用することができない時は、コンピュータ・センターの重役秘書を買収する。秘書はたいして貴重な情報に簡単にアクセスすることができ、しかも安い賃金で働いているはずだ。人的要因の効果をめざってはいけな。

これ以外にも各種の攻撃方法がLinde(1975 NCC)によって紹介されている。

5.4.4 セキュリティの設計理念

Saltzer および Schroeder(1975)は、安全なシステムを設計するためのいくつかの基本理念を発表した。両者の理念を以下に簡単にまとめる(MULTICS での経験に基づいている)。

第1に、システムの設計技法を公開する。侵入者がシステムの仕組みを知らないだろうと、過信してはならない。

第2に、デフォルト値をアクセス不可にする。合法的アクセスが拒絶されるようなエラーは、不正アクセスが検出される前に早く報告する。

第3に現在与えられている権利を調べる。システムは、許可を調べ、アクセスを許すかどうかを決定するが、その情報を後で使用するために保存してはならない。多くのシステムはファイルを開く時に許可内容を調べ、以後確認は行っていない。したがってユーザーがファイルを開いたまま何週間も経過すると、すでにファイルの所有者がファイルの保護内容を変更してかなり経過していても、ユーザーは変更前の保護内容のままでそのファイルに引き続きアクセスすることができ。

第4に各プロセスに最低限の特権しか与えないようにする。エディタが特定のファイル(エディタを呼び出す時に指定した)を編集するためのアクセス権しか持っていないとしたら、トロイの木馬の手法を用いたエディタもあまり効果はない。この理念はきめ細かい保護手法を示している。この章の後半でその様な手法に関して説明する。

第5に保護機構は単純で、一貫性を持っており、システムの最下位に組み込まれていない限りならぬ。保護機能の弱い既存システムに合わせてシステムを改造することはできない。保護機構は、正当性と同じで後から追加できる機能ではない。

第6に選んだ手法は心理的に受け入れやすいものでなくてはならない。ユーザーがファイルの保護を面倒だと感じるなら、保護は行われないだろう。しかしいったん問題が発生すると、不満を爆発させることが多い。それはあなたのせいです。* といってもそれは受け入れられない。

5.1 セキュリティ

5.4.5 ユーザーの認証

多くの保護手法は、システムがそれぞれそのユーザーを識別できるという仮定に基づいている。ユーザーがログインした時に実行される識別は、認証(user authentication)と呼ばれている。多くの認証方法はユーザーの知っていること、持っているもの、ユーザー自身などに基づいている。

■ パスワード

認証の最も一般的な形態は、ユーザーにパスワードを入力させることである。パスワードを使った保護はわかりやすく、実行も簡単である。UNIX では以下の様な動きをする。ログイン・プロセスがユーザーに名前とパスワードを入力するよう促す。パスワードは直ちに暗号化される。次にログイン・プログラムは1ユーザーあたり1行のASCII文字列から構成されるパスワード・ファイルを読み取り、ユーザーのログイン名を含む行を抽出する。この行に含まれたパスワード(暗号化された)が、入力したものを暗号化したパスワードと一致すれば、ログインは認められ、そうでなければログインは失敗に終わる。

パスワードによる認証を破るのも簡単である。事実、大企業や政府機関の所有する機密システムに家庭用コンピュータを使って侵入した高校生や、中学生グループに関する記事を頻りに見かける。本質的には侵入は、いつでもユーザー名とパスワードの組み合わせを予測することによって行われる。

Morris および Thompson(1979)は、UNIX システムにおけるパスワードの研究を行った。彼らは一般的なパスワードをリストにまとめた。リストによると、姓名、氏名、通りの名前、町名、普通の辞書に収録されている言葉(そしてそれを逆に読んだもの)、車のナンバープレートの番号、ランダムな文字から成る短い文字列などが用いられやすいという。

そして次にこれらを1つずつ既知のパスワード暗号化アルゴリズムを用いて、暗号化した。そして暗号化されたパスワードのうちリストのエントリと一致するものがないか調べたところ、パスワードのうち86%がリストに記載のものであったという。

印字可能な95のASCII文字の中からランダムに抽出された7文字だけをパスワードに使用する場合、検索空間は95⁷、すなわちおよそ7×10¹²となる。1秒あたり1,000回暗号化作業を行う場合、パスワード・ファイルとの照合を行うためのリストを作成するには、2,000年もかかることになる。さらにリストは磁気テープ20,000,000本もの長さになる。パスワードに小文字、大文字、そして特殊文字をそれぞれ少なくとも1つ以上含ませ、7文字以上にするだけでユーザーが選択するパスワードはかなりの改善をもたらす。

ユーザーに適切なパスワードを選択させることが実質的に不可能なとしても、Morris と Thompson はその様な攻撃(あらかじめ多数のパスワードを暗号化しておく)を無意味にする技法を紹介している。彼らのアイデアはnビットのランダムな番号を、各パスワードに連結するものである。パスワードが変更されればランダムな番号も変更される。ランダムな番号は暗

号化される前の形でパスワード・ファイルに保存されており、誰にでも読み出せるようになってい
る。暗号化されたパスワードをパスワード・ファイルに保存するため、パスワードとランダムな番
号が最初に連結され、共に暗号化される。暗号化されたものはパスワード・ファイルに保存される。
ここでパスワードに使用されるランダムな番号を作った暗号化し、その結果を分類済みのファ
イルに保存し、暗号化された任意のパスワードを簡単に検索しようとしている侵入者を考えて
みよう。侵入者がパスワードがMarilynであると子思したら、Marilynを暗号化してその結果を
1に入れるだけでは不十分である。Marilyn0000, Marilyn0001, Marilyn0002 など、26分の暗号化
を行い、すべてを1に保存しなくてはならない。この技法によって1の大きさも2になる。UNIX
はこの技法を $n=12$ で使用している。

この技法はあらかじめ暗号化された多数のパスワード・リストを作成することによって侵入を
たくらむ者に封じては保護効果があるが、David というユーザーが、そのパスワードも David と
している場合には、なんの効果もない。より適切なパスワードを選択させるためには、コンピュ
タになんらかの提案を行わせるとよい。コンピュタによっては forally, garbunzy または
hippity などパスワードとして使用できる。発音の簡単な、しかも意味をなさぬ言葉(でさるだけ
大文字と特殊文字が盛り込まれているもの)を作成するプログラムを持っているものもある。

また別のコンピュタではユーザーにパスワードを定期的に変更させ、パスワードが誰かに知
られた場合、被害を最小にとどめている。この技法の極端な例は一時パスワード(one time
password)である。一時パスワードが使用される場合、パスワードのリストを含んだ冊子を渡され
ることになる。ログインのたびにリスト内の次のパスワードを用いる。侵入者がもしパスワード
を見つけても、次回から異なるパスワードを使用することになるので、どうにもならない。ユー
ザーはこの冊子をなくさないように心がけなくてはならない。

いまでもなく、パスワードが入力されてもコンピュタは入力された文字を表示してはなら
ない。端末の近くに悪意のユーザーがいる可能性があるからである。またパスワードを暗号化さ
れる前の形態でコンピュタに保存しないこと、そしてコンピュタ・センターの管理側にも暗号化
される前の形態で保存しないことなども重要である。暗号化される前のパスワードがどこに保存
されても、問題になりやすい。

また新しいユーザーに長い質疑応答リストを提供させ、暗号化された形態でそれをコンピュ
タに保存するといったパスワード設定技法もある。質問内容はユーザーが書き留めなくても選択
できるような形式のものである。典型的な質問例を以下に挙げる。

- Marjolein の妹の名は?
- 小学校が位置していた通りの名前は何?
- Mrs. Waroboff が教えている教員は?

ログイン時にコンピュタはランダムに1つの質問を抽出し、その答を確認する。
もう1つの方法はチャレンジ・レスポンス(challenge-response)と呼ばれている。この技法を用

いると、ユーザーはユーザー登録中に例えばアルゴリズムとして 2^x を選択する。そしてユーザ
がログインすると、コンピュタは引数として例えば7を表示し、その場合ユーザーは49と入力
する。アルゴリズムは道の幅や、端末、また朝と夜でも異なっている可能性がある。

物理的識別

全く異なった設置方法として、ユーザーがなんらか、例えば磁気ストライプを貼っているプラ
スチック・カードを持っているかどうかを確認するというやり方もある。カードを端末に挿入さ
せ、そこで誰のカードであるかを判別する。この方法をパスワードと組み合わせて、(1)カードを
持っており、(2)パスワードを知っているユーザーのみログインを許可することもできる。自動預
金引出機がこの典型的な例である。

さらにもう1つの方法は、製造の困難な物理的特徴を測定することである。例えば端末に指紋
または声紋認識装置を付けてユーザーを識別することができ(ユーザーがコンピュタに自分が
誰であるかを告げ、特定の指紋と照合すれば、データベース全体を照合するよりも検索は速く行
える)。直接目で見て判断することはまだできないが、今後その様な機能が提供される可能性もあ
る。

そしてもう1つのアプローチとして署名を用いるものもある。ユーザーは端末に登録された特
殊なペンで署名し、コンピュタがそれを意識されている署名と比較する。署名の上手下手では
なく、署名中のペンの動きと比較する。署名を真似ることはできても、どのストロークを先に打
つかまで知ることはできない。

これ以外にも指の圧力を分析するという方法があるが、これは非常に効果的である。この技法
を用いるためには、各端末に図5.21の様な装置を用意する。ユーザーは手をそこに挿入し、すべ
りの指の圧力が測定され、データベースと照合される。

さらに例をあげ続けることもできるが、以下の2つを重要例として紹介する。ネコやその他の
動物は自分の領域の周りに徘徊することによって、その境界を守ろうとする。例えば使役用の器
械を用いて非常に簡単な識別を行なったとしよう。各端末にはこの機械が1台備わっており、それ
ぞれ「ログインの際はサンプルをここにに入れてください」という内容のメッセージが現れる。こ
れは絶対投入することのできないシステムだが、おそらくかなり深刻なユーザー側の受け入れに
関する問題が発生するだろう。

同様なことが面紙と小型の分光器から成るシステムにもいえる。ユーザーは面紙で指紋を転す
よう指示され、それによって出た血を分光器による分析にかけ、重要な点は、どの様な設置方
法であれユーザー社会において精神的に受け入れられなければならないということである。指の
長さや刻む程度なら特に問題はないが、コンピュタに指紋を保存しておくことすら受け入れた
くない人も多い。

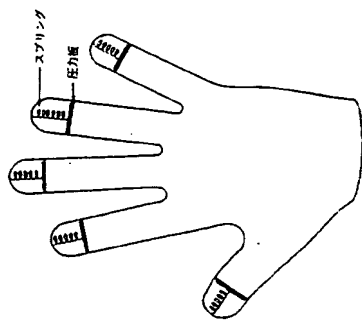


図 5.21 指の長さで測定する装置

■ 対処方法

侵入者がシステムに入り込んで、多大な被害を引き起こした後、嚴重な保護機構を採用することになったシステムでは、不正な入力に非常に困難になるような措置を取ることが多い。例えば各ユーザーは特定の端末からしかログインすることができず、しかもログインできるのは特定の曜日と、特定の時間といった制限付きであることも多い。

ダイヤル回線を以下の様に使用するのも1つの方法であらう。それでもダイヤルを回して、ログインすることはできるが、ログインが成功すると、システムは直ちに回線を切り、指定された番号のユーザーを呼び返す(コールバック)。これによって侵入者は仕様のダイヤル回線から侵入することができなくなり、ユーザーの(特定の)電話だけにその権利が与えられることになる。コールバックを行う場合も、行わない場合もシステムはダイヤル回線から入力されたパスワードを検査するために少なくとも10秒を要する。ただしこの時間は回線が連続してログインに失敗すると延長され、侵入者の侵入率を低下させることができる。また、ログインに3回失敗すると、回線は10秒間断線され、その間に保護管理者に通報するというのも効果的である。

ログインの記録はすべて保存しておくなければならない。ユーザーがログインすると、システムは前回のログイン時間と端末を報告し、それにより侵入を知ることができる。簡単なパスワードを持つ特殊なログイン名(例えばログイン名: guest, パスワード: guest)を提供する方法が簡単である。この名を用いてログインしたものがあればシステム保護管理者に通報される。これ以外の方法として、オペレーティング・システムに発見しやすいバグを仕掛けたり、侵入者の権限を目的とした同様の仕掛けを行うことである。

5.5 保護機構

ここまでの節では、多くの問題点について見てきた。その中には、技術的なものもあれば、そうでないものもあった。この節ではファイルその他の保護を行うためにオペレーティング・システムで用いられているさまざまな技法を解説していきたい。これらの技法はすべて方針(どのデータを誰から守るのか)および機構(方針に従うためにどのようなシステム形態をとるか)の区別が明確になされている。方針と機構の区別に関しては Levin 他(1975)を参照されたい。ここでは方針ではなく、機構に重点をおいて説明していく。

5.5.1 保護ドメイン

コンピュータ・システムには保護を必要とする対象物(objects)が多数含まれている。対象物は、CPU、メモリサブメント、端末、ディスクドライブ、プリンタなど、ハードウェアである場合もある。プロセス、ファイル、データベース、セマフォといったソフトウェアである場合もある。各対象物には、参照に用いられる一意な名前が付けられており、また対象物に対して実行可能な一連の操作を持っている。ファイルに対しては READ および WRITE が、またセマフォには UP と DOWN が適切である。オペレーティング・システムにおける対象物は、プログラミング言語での、抽象データ型(abstract data types)と呼ばれるものに相当する。

いうまでもなく、プロセスがアクセス権を与えられていない対象物にアクセスできないようにするための方法が必要である。さらにこの機構は、必要に応じてプロセスによる特定の正当な操作も禁止できなくてはならない。例えばプロセス A はファイル F の読み取りは行えても、書き込みはできない、といった指定も可能である。

異なる保護機構について述べるために、ドメインに関する説明を加えておくことにする。ドメイン(domain)は(対象物、権利)のペアの集合である。各ペアは対象物とそれに対して行うことのできる操作の部分集合を構成している。ここでの権利(rights)とは、これらの操作のいずれかを実行するために与えられる権利を指す。

図 5.22 には 3 つのドメインを示しており、各ドメインには対象物とその権利(Read, Write, execute)が含まれている。プリンタが同時に 2 つのドメインに存在している点に注意されたい。この例では示していないが、同じ対象物を複数のドメインに置き、各ドメインに異なる権利を与えることも可能である。

プロセスは、常になんらかの保護ドメイン内で実行されている。すなわちアクセス可能な対象物の集まりが存在し、各対象物に対してなんらかの権利を与えられている。プロセスは実行中に 1 つのドメインから別のドメインに移動することも可能である。ドメインの切り換えに関する規定はシステムによって大きく異なっている。

保護ドメインの概念をさらに明らかにするための、UNIX を見てみよう。UNIX ではプロセスの

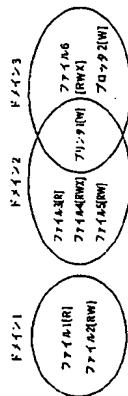


図 5.22 3つの保護ドメイン

ドメインは uid と gid で決定される。どの様な (uid, gid) の組み合わせを用いても、アクセス可能な対象物(特殊ファイル)として表現されている出力装置も含めたファイル)をすべてリスト表示したり、読み書き実行のいずれに對するアクセス権を与えられているのかを知ることができる。同じ (uid, gid) の組み合わせを持つ2つのプロセスは、同一の対象物の集合にアクセスできる。異なる (uid, gid) 値を持つプロセスは、異なるファイルにアクセスを行うことになるが、たいていかなりの量のファイルが重複している。

さらに UNIX のプロセスはユーザー部分とカーネル部分を半分ずつ持つっており(3章の図3.14)、カーネル部分はユーザー部分と異なる対象物のセットにアクセスを行うことができる。例えばカーネルは物理メモリ内のページ、ディスク全体、その他保護されたすべてのアクセスすることができる。したがってシステムコールはドメインの切り換えを引き起こす。

プロセスが SETUID または SETGID ビットを持つファイルを実行した場合、新しい有効 uid または gid を得る。異なる (uid, gid) の組み合わせを用いると、異なるファイルと操作セットを指定することができる。SETUID または SETGID を持つプログラムの実行もドメインを切り換えることになる。

UNIX では、プロセスをカーネル部分とユーザー部分に分割しているが、これは MULTICS で用いられていた、より強力なドメイン切り換え機構を受け継いだものである。MULTICS では、ハードウェアがサポートしていたプロセスあたりのドメインは2つ(カーネルとユーザー)ではなく、最大64であった。MULTICS のプロセスはそれぞれがなんらかのドメインで実行中の手続きから構成されていた。ここのドメインはリング(rings)と呼ばれていた(Schroeder, Saltzer, 1972)。手続きは実行中のプロセスに動的にリンクされることである。

図 5.23 は4つのリングを示したものである。最も内側のリングはオペレーティング・システムのカーネルであり、最も強力である。カーネルから外に移動するにつれ、リングの力は弱くなる。例えばリング1には UNIX における mkdir などの様な、root が所有し、SETUID を持つプログラムによって処理される機能が含まれている。リング2は学生のプロセスを評価するためのプログラムが、そしてリング3には学生プログラムが含まれている。

1つのリング内の手続きが別のリングの手続きに呼び出されると、トラップが発生し、システムはプロセスの保護ドメインを変更することができる。したがって MULTICS プロセスは、一生のうち最大64の異なるドメインで実行されることができた(実際には上述の状況よりも複雑なこ

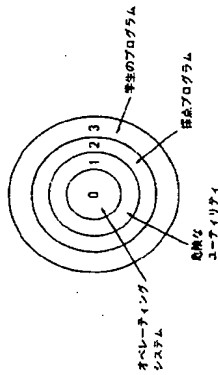


図 5.23 4つのリングを用いた MULTICS のプロセス。各リングは異なる保護ドメインに存在する。

とが多い。手続きは複数の隣接するリングで実行することが可能で、リング間でやりとりされる引数は入念に制御されていた。MULTICS に関する詳細は、Organick (1972) を参照されたい。それでは、システムはどの様な方法でどの対象物がどのドメインに属しているかを記録しているのだろうか。概念的には少なくともドメインを行とし、対象物を列とする大規模なマトリックスを想像することができる。図 5.24 のマトリックスを示したのが図 5.24 である。このマトリックスと現在のドメイン番号を用いて、システムは常に指定したドメインから特定の方法で特定の対象物に対するアクセスが可能かを判断できる。

対象物

ドメイン	ファイル1	ファイル2	ファイル3	ファイル4	ファイル5	ファイル6	プリンタ	プロセス
1	Read	Write						
2			Read	Write	Read	Write	Write	
3							Read	Write

図 5.24 保護マトリックス

ドメインの切り換え自体は、MULTICS と同様、ドメイン自身が対象物であるとし、ENTER の操作を行うことによって、マトリックス・モデルに容易に含ませることができる。図 5.25 は図 5.24 のマトリックスを再度図示したもので、今回は3つのドメインを対象物自身として表している。ドメイン1のプロセスはドメイン2に移行することができるが、移行した後はもうもとに戻ることはできない。この様な状態は、UNIX における SETUID を持つプログラムの実行を表したものである。この例ではこれ以外のドメイン切り換えは認められない。

5 章 ファイルシステム

5.5 保護機構

図 5.25 対象物としてのドメインを持つ保護マトリックス

ドメイン	File0	File1	File2	File3	File4	File5	File6	File7	File8	File9	File10	File11	File12	File13
1	Read	Write												Enter
2			Read	Write	Read	Write								
3							Read	Write	Execute					

図 5.25 対象物としてのドメインを持つ保護マトリックス

5.5.2 アクセス制御リスト

実際には図 5.25 のマトリックスは非常に大きく、分散しているため、これを保存することは滅多にない。大半のドメインはほとんどどの対象物に対して全くアクセス権を与えていないため、大きくて、中身の無いマトリックスを保存することはディスク空間の浪費につながる。それでもマトリックスの保存が必要な場合は、行または列単位で、しかも空きエレメントは除いて保存するようにする。2つのアプローチは、互いに異なっているため、この項では列単位での保存について触れ、次の項において行単位での保存について触れることにする。

最初の技法では、各対象物にリスト(一定の順番を付けた)を連結させる。このリストには対象物にアクセスすることのできるすべてのドメインと、その方法が記載されている。このリストはアクセス制御リスト(access control list)または ACL と呼ばれている。UNIX で実現する場合、最も簡単な方法は、ACL を各ファイルの個別ディレクトリブロックに置き、ファイルのモードにブロック番号を記録することである。マトリックスのうち、空のものは対象となっていないため、すべての ACL 保存容量はマトリックス全体に比べるとかなり少なくて済む。

ACL の仕組みを知るために、引き続き UNIX において(uid, gid)ペアでドメインが表現される例を考えてみよう。事実 ACL は UNIX の前身である MULTICS で、以下の様な形態で使われていた。したがって、ここで挙げる例は単なる空想上のものではない。

ここでは Jan, Els, Jelle, Maaike という 4 人のユーザー(すなわち uid)があると仮定しよう。それぞれのユーザーはシステム、スタッフ、学生、学生というグループに属している。いくつかのファイルが以下の ACL を持っているとする。

```
File0: (Jan, *, RWX)
File1: (Jan, system, RWX)
File2: (Jan, *, RW-), (Els, staff, R-), (Maaike, *, R-),
File3: (*, student, R-),
File4: (Jelle, *, -), (*, student, R-)
```

カッコ内の ACL エントリは uid, gid, および許可されているアクセス内容(Read, Write, eXecute)を示している。アスタリスクはすべての uid または gid を表している。File0 は uid=Jan なら、任意の gid を持つプロセスによって読み書き実行を行うことができる。File1 は uid=Jan で gid=system のプロセスだけがアクセス権を与えられている。uid=Jan で gid=staff のプロセスは File0 にアクセスすることはできても、File1 にアクセスすることはできない。File2 は uid=Jan で、任意の gid を持つプロセスか、uid=Els で、gid=staff のプロセス、または uid=Maaike で、任意の gid を持つプロセスのいずれかによって読み書きが行える。File3 は学生ならば誰でも読取ることができる。File4 は特に興味深い。uid=Jelle ならば、グループに関連なくアクセスを全く認めないものである。ただし他の学生はそれを読取ることができる。ACL を用いると、同じクラスの他の者にはアクセスを与えながら、特定の uid または gid による対象物のアクセスだけを禁じることができる。

この様な機能は UNIX にはない。それでは UNIX に提供されているものを見てみることにしよう。UNIX ではファイルの所有者、所有者のグループ、一般ユーザーそれぞれに対して 3 つのビット rwx を提供している。手法は ACL と同じだが、9 ビットに縮められている。これは対象物に連結されたリストで、誰がどの様にその対象物にアクセスできるかを示している。9 ビットを用いた UNIX のスキーマは ACL システムに比べてあまり一般性はないが、実際にはこの技法の方が優れており、実現も簡単でありコストも安い。

対象物の所有者はいつでも ACL を変更し、それまで与えられていたアクセス権を簡単に禁止することができる。唯一の問題となるのは、ACL を変更しても現在その対象物を使用しているユーザー(例えばファイルをオープンしているなど)には、おそらくなんの影響もないという点である。

5.5.3 権限

図 5.25 のマトリックスを行単位で切り分けることもできる。この方法を用いると、各プロセスにはアクセスすることのできる対象物のリストと、それぞれに認められている操作、すなわちドメインが連結される。このリストは権限リスト(capability list)と呼ばれており、そこに置かれた個々の項目を権限(capabilities)と呼んでいる(Dennis Van Horn, 1966; Fabry, 1974)。

典型的な権限リストを図 5.26 に示している。それぞれは対象物の権限を示すフィールド type と、この種の対象物に認められた正当な操作を表したビットマップであるフィールド rights。そして対象物自身を指すポインタであるフィールド object(例えばノード番号など)を持つ。権限リスト自体も対象物であるため、他の権限リストからポインタを使って指すことも可能である。これはサブドメインの共有を可能にする。許可は権限リストにおけるその位置によって参照されることが多い。プロセスが権限 2 が指しているファイルから 1K を読み出せ* という要求を行ったとしよう。この様なアドレス付けの形態は、UNIX におけるファイル記述子の使用と似ている。権限リスト、すなわち C リスト(C-lists)は、ユーザーによる不正な変更から守るのではなく

タイプ	Rights	対象物
0	Read	ファイルのポインタ
1	Write	ファイルのポインタ
2	Execute	ファイルのポインタ
3	Write-Execute	ファイルのポインタ

図 5.26 図 5.24 におけるドメイン 2 の権限リスト

ない。そのために3つの方法が提供されている。最初の方法ではタグアーキテクチャ(tagged architecture)と呼ばれるものが必要である。このアーキテクチャでは、権限の有無を示す子欄(またはタグ)ビットを各メモリワードに付加する。タグビットは演算、比較、その他通常の命令には使用されない。そしてカーネルモードで実行中のプログラム(すなわちオペレーティングシステム)だけがそれを変更することができる。

2番目の方法としてはオペレーティングシステムの内部にCリストを保存し、上記の様にスロット番号を使ってプロセスが参照を行うというものである。Hydra(Wulf, 1974)はこの様な構造を用いていた。3番目の方法としては、Cリストをユーザー空間に置きながら、各権限をユーザーの知らない秘密のキーで暗号化するのである。この方法は特に分散システムには効果的であり、Amoebaでも使用されている(Tanenbaum 他, 1986)。

4つ目や実行など対象物に依存する権利のほかに、すべての対象物に適用される一般的な権利(generic rights)が権限に提供されている。一般的権利の一例を以下に挙げる。

- 権限コピー: 同一対象物に新しい権限を作成する。
- 対象物コピー: 新しい権限を持つ複写対象物を作成する。
- 権限削除: 対象物とそのままで、Cリストからエントリを削除する。
- 対象物破壊: 対象物と権限を永久に消去する。

多くの権限システムは、対象物の各タイプに対してタイプマネージャ・モジュール(type manager module)を持つ。モジュール集合から構成されている。ファイルに対して操作実行を促す要求はファイルマネージャに、そしてメールボックスに関する要求はメールボックス・マネージャに送られる。これらの要求は該当する権限を伴っている。タイプマネージャ・モジュールは通常のプログラムであるため、ここで問題が生じる。ファイルの所有者はファイルの操作のいくつかを実行することができるが、その内部表現(例えばiノードなど)を得ることはできない。タイプマネージャ・モジュールに通常のプロセス以上の権限を持たせることが必要である。

Hydraでは権利の拡大(right amplification)という技法でこの問題を解決した。この技法ではタイプマネージャに特定の権利テンプレートを与え、このテンプレートによってそれ自身が認められているより多くの権利をオブジェクトに与えることができる。対象物を強く型付けし、分類している他の権限システムでも、同様の技法が用いられている。

権限システムに関し最後に触れなければならないのが、対象物に対するアクセスを無

効にすることがかなり困難であるという事実である。システムが任意の対象物に対して外部にある権限を授け出し、操作することは困難である。ディスク上に分散されているCリストに保存されている可能性もあるからである。1つの解決策として各権限に対象物自身ではなく、間接対象物を指定するという方法がある。間接対象物、実際の対象物を指定することにより、システムは常にその境界を断ち、権限を無効にすることができる(後から間接対象物に対する権限がシステムに提示されると、ユーザーは間接対象物がなにも対象物を指していないことを発見する)。

許可を無効にするためのもう1つの方法として、Amoebaで使用されている手法がある。各対象物は長いランダムな番号を含んでおり、これと同じものが権限にも提供される。権限を使用するために番号を提示すると、その2つの番号が比較される。番号が一致した場合に限り、操作が認められる。対象物の所有者は、対象物のランダムな番号の変更を要求し、既存の権限は無効にすることができる。これら手法のいずれも選択的な無効。つまりJohnには権限を与え、他者には与えないといった指定を行うことはできない。

5.5.4 保護モデル

図 5.24 の保護マトリックスは静的なものではない。新しい対象物の作成、古い対象物の破壊、そして所有者が対象物に対するユーザー一般の増減を行う際に、頻繁に変更が行われる。保護システムの構築はかなり注意深く行われており、それにより保護マトリックスは絶えず変更されている。以下の項ではこの作業について簡単に触れることにする。

Harrisonら(1976)は、あらゆる保護システムの基本として用いることのできる、保護マトリックスに対する6つの基本操作を明確化した。これらの操作とはCREATE OBJECT、DELETE OBJECT、CREATE DOMAIN、DELETE DOMAIN、INSERT RIGHT、REMOVE RIGHTである。最後の2つの基本操作は特定のマトリックス・エレメントに権利を挿入したり、削除したりするものである。その一例としてドメイン1にFile6を譲取る権利を与えることが挙げられる。

これら6つの基本操作は、保護コマンド(protect commands)として組み合わせることでもできる。ユーザープログラムは保護コマンドを使用してマトリックスの変更が行える。この場合、基本操作を直接実行することはない。例えばシステムは新しいファイルの作成をコマンドで行い、その際ファイルが既存のものでなければ、新しい対象物を作成し、所有者にすべての権利を与えるという作業を行うこともある。また所有者が、システムに存在する各ユーザーに対して譲取りの権利を与えるようなコマンドの存在も考えられる。このコマンドによって各ドメインの新しいファイルのエントリに「譲取り」の権利を挿入することになる。

マトリックスは常に、認証ではなく、そのプロセスがどのドメインに属しているかによって権利を決定することになっている。マトリックスはシステムがどの様に権利を行うかを決定し、認証はどの様に管理するかを決定している。この2つの違いを明らかにしている例として、図 5.27 の簡単なシステムを考えてみよう。ここではドメインがユーザーに相当している(UNIX モデル

5章 ファイルシステム

に類似している。図5.27(a)では基図した保護方針が達成されている。すなわち Henry は mailbox7 を読み書きでき、Robert は secret の読み書きができ、3人のユーザーはすべて compiler の読み書きが可能である。

対象物		対象物	
compiler mailbox7 secret		compiler mailbox7 secret	
Read	Execute	Read	Execute
Write	Execute	Write	Execute
Read	Write	Read	Write
Write	Write	Write	Write
Read	Read	Read	Read
Write	Read	Write	Read
Read	Write	Read	Write
Write	Write	Write	Write

図5.27 保護マトリックス

(a) 隠蔽されている状態 (b) 隠蔽されていない状態

それは Robert が非常に頭がよく、マトリックスを図5.27(b)の様に変更するコマンドを発見したとして、ここでアクセス権のない mailbox7 に対するアクセスを行った。読出しを行おうとするとオペレーティングシステムは図5.27(b)の状態が認められていないことを知らないうちにその要求を実行してしまう。

使用可能なマトリックスの集合を2つの集合に区分できることがわかった。すなわち隠蔽されている状態と、隠蔽されていない状態の2つの集合である。多くの理論的な研究がもたらした疑問は、「最初の隠蔽状態から、一連のコマンドを与えられても、システムが隠蔽されていない状態にならないことを証明できるか」というものである。

すなわち与えられた隠蔽(保護)コマンドが特定の保護方針の強化に通じているかという疑問である。方針の簡単な例として、軍事に用いられるセキュリティ手法を考えてみよう。各対象物は、秘密ではない、秘密、秘密、または極秘扱いのいずれかである。各ドメイン(したがって各プロセスでもある)も、これら4つの秘密保護レベルのいずれかに属する。セキュリティの方針では、以下の2つの決まりがある。

- プロセスが自分より高いレベルの対象物を読取することはできないが、同等あるいはそれ以下のレベルの対象物は自由に読取ることができる。秘密レベルのプロセスは秘密レベルの対象物を読取することはできないが、極秘レベルの対象物を読み取る権利はない。
- プロセスは自分より低いレベルの対象物に情報を書き込むことはできない。秘密レベルのプロセスは極秘レベルのファイルに書き込むことはできるが、極秘レベルのファイルには書き込めない。

5.5 保護機構

軍事用語に置き換えて、兵士が戦線、下士官が戦線、将校が戦線レベルで活動を行っている。仮定すると、下士官は兵士の書類を見ることはできても、将校の書類を勝手に見ることはできない。同様に下士官は自分の知っていることを将校に報告することはできても、兵士はその権利を持たないため、彼らに報告してはならない。

この様な方針から、あるマトリックスの初期状態(各対象物のレベル決定方法も含む)から、一連のコマンドによりマトリックスが修正されていく過程で、システムの安全性を規定するにはどうしたら良いのであろうか。これを完全に言うことは非常に困難であることがわかった。多くの汎用システムは理論的に言うに安全性に欠けている。この点に関する詳細は Landwehr (1981) と Denning (1982) の報告書を参照されたい。

5.5.5 隠れたチャネル

以上の項では、保護システムに対する正式なモデルの作成方法を見てきた。しかしここではその様なモデルを作成することがどれだけ無益であるかを学ぶことになる。特に安全性が絶対的に保証されているシステムにおいては、理論的には通称が不可能なプロセス間の情報漏洩は比較的簡単に行われている。これに関しては、Lampson (1973) の報告を参照されたい。

Lampson のモデルには3つのプロセスが用いられており、主として大型のタイムシェアリングシステムに適用されている。最初のプロセスはクライアントで、第2プロセスであるサーバーになんらかの操作を要求するものである。クライアントとサーバーは互いに安全に通信している。例えば、サーバーの動きがクライアントの租税申告書作成を手伝うことであるとする。クライアントは、サーバーが秘密裏に、例えば個人の収入リストを保存して、リストを誰かに売り付けたりしないかを心配している。サーバーはクライアントがこの重要な租税プログラムを盗まれないかを心配している。

第3のプロセスは敵側のスパイで、サーバーにクライアントのデータを盗取らせよう入れ知恵する。スパイとサーバーは通常同じ所有者のものである。図5.28では3つのプロセスが示されている。この計画の目的はサーバーが合法的な手段でクライアントから入手した情報をスパイに漏らすことができないようなシステムを構築することである。Lampson はこれを拘束問題 (containment problem) と呼んだ。

システム設計者の見方からいえば、サーバーがスパイに情報を漏らせないように、サーバーのオブジェクト化、または拘束を行えばよい。保護マトリックス手法を用いれば、スパイが読み出せるファイルにサーバーが書き込みを行うことで、サーバーとスパイが通信することを防止することができる。またシステムのプロセス間通信機構を用いる場合も、同様なことを保証できるだろう。しかし残念ながら、より狡猾な通信経路が存在している。例えばサーバーは以下の様に2通りのビット列を送信することができる。1ビットを送るためには、特定の時間間隔の間で限りの演算処理を行う。0ビットを送る場合には、同じ時間の間、スリープ状態となる。

スパイはその応答時間を監視することによって、ビット列を検出する。一般的には、サーバー

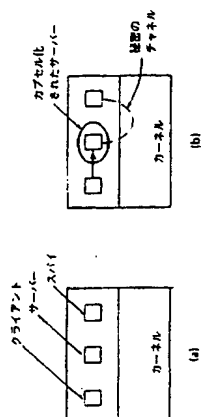


図 5.28 隠れたチャネルによる情報漏洩
(a)クライアント、サーバーおよびネットワークの接続
(b)キャッシュとコアの接続による情報漏洩

が1を送信している時よりも、0を送信している時の方が応答はよくなる。この通信経路は隠れたチャネル(covert channel)として知られている。図 5.28(b)はこれを図示したものである。

いまでもなく、隠れたチャネルは多くの重要な情報を含むための雑音の多いチャネルではあるが、エラー訂正コード(例えばハミングコード)あるいはより高度なもの(例えば、雑音の多いチャネルでも確実に情報を送ることができる、エラー訂正コードを使用すると、すでに低い帯域幅をさらに低めることになるが、それでも重要な情報を漏洩することは不可能である。対象物とドメインのマトリックスに基づいた保護モデルを使って、この種の漏洩を防ぐことはできない。

CPUの利用率を使って情報を漏洩することだけが隠れたチャネルではない。ページング率も実際に利用することができる。ページングファクトリが存在すれば、存在しなければ、畢竟、何らかの方法でシステム性能を低下させるのであれば、実際に利用できる可能性がある。システムがファイルのロック機能を提供している場合、サーバーはファイルをロックして1と表示し、ロックを解除したら0と表示することができる。通常、アクセスすることのできないファイルに属してもロック状態を知ることができる。

占有リソース(テープドライブ、プロッタなどの入手と解除も同様)に利用することができる。サーバーは資源を獲得して1を送信し、それを解除して0を送信する。UNIXではサーバーがファイルを作成して1を表示させ、それを削除して0を表示する。スバビはファイルの存在を確認するためにシステムコール ACCESS を用いることもできる。このコールはスバビがファイルを使用する権利を持っていないことも通ずる。残念ながらこれ以外にも多くの隠れたチャネルが存在する。

Lampson は、さらにサーバープロセスの所有者に情報を漏らす方法についても触れている。サーバープロセスは、クライアントに料金を請求するために、所有者にクライアントのためにどれだけの作業を行ったかを報告する。実際の請求料金が例えば 100 ドルで、クライアントの収入が 53,000 ドルであれば、サーバーは請求書に 100.53 ドルと書き込み、所有者に報告できる。

すべての隠れたチャネルを捜し出し、防御を行うことだけでも非常に困難である。実際にはそ

れをどうすることもできない。ランダムにページフォルトが発生するプロセスを除外したり、隠れたチャネルの帯域幅を低下させるためにシステム性能を下げることはあまり望ましくない。

5.6 MINIX ファイルシステムの概要

他のファイルシステムと同様に、MINIX のファイルシステムもこれまでで学んできた問題すべてに直面しなければならぬ。ファイルのための空間を割り当てたり、削除したり、ディスタックロックや空き空間の記録を取ったり、不正な使用からファイルを保護したり、といった問題が山積みされている。以降では MINIX がこの問題をどの様な方法で解決しているかを詳しく解説していく。

本章では、前節までは一般性を重んじ、MINIX ではなく UNIX を何度となく説明に用いた。もちろんこの2つの外部インターフェイスは全く同じである。ここでは MINIX の内部構造について見ていくことにしよう。UNIX の内部構造に関しては、Thompson(1978)と Bach(1986)の記述を参照されたい。MINIX のファイルシステムは単なるユーザー空間で動作する C のプログラムの要約(4章の図 4.20 参照)。ファイルを読み書きする際は、ユーザープロセスがメッセージをファイルシステムに送信し、必要な作業を要求する。ファイルシステムは指定された作業を行い、返答する。ファイルシステムは実質上、呼出し側と同じマシンの上で動いているネットワーク・ファイルサーバーにすぎない。

この様な構造は重要な意味を持っている。その1つは、ファイルシステムを MINIX の他の部分とは独立して修正、実験、検査できるようにすることである。また、C コンパイラを持つコンピュータであれば、ファイルシステム全体を移植し、そのコンピュータ上でコンパイルし、UNIX に類似した独立型のリモート・ファイルサーバーとして使用することもできる。この時変更しなければならぬのは、システムによって異なるメッセージの送受信方法に関する記述だけである。

以下の項では、ファイルシステムの設計においてキーポイントとなる項目に關し、その概要を述べる。なかでもメッセージ、ファイルシステムの構造、ビットマップ、i ノード、ブロックキャッシュ、ディレクトリ、パス名、プロセステーブル、特殊ファイル(およびパイプ)に關して詳しく見ていくことにする。これらに關する説明の後、ユーザープロセスがシステムコール READ を実行した時のそれぞれの状態を追跡し、各部分ごとの様に連携しているのかを調べる。

5.6.1 メッセージ

ファイルシステムは、作業の要求を行う 29 種類のメッセージを受け付ける。そのうち 2 つを除いた残りすべてが、MINIX システムコールに對するものである。2 つの例外は、MINIX 内のファイルシステム以外の部分で生成されるメッセージである。すべてのメッセージと、その引数、および結果を図 5.29 に示している。ファイルシステムは、メモリアドレスからも、メモリアネー

iノードは、さらにモード情報も含んでおり、ファイルの種類(通常ファイル、ディレクトリ、ブロック型特殊、キャラクター型特殊、あるいはパイプ)、保護状態とSETUID、SETGIDビットを提供する。iノードレコード内のフィールドであるリンク数は、iノードを指しているディレクトリエントとiノードの数を表している。ファイルシステムは、このリンク数により、ファイルに割り当てられている領域を解放するかどうかを判断している。このフィールドを、ファイルの現在のオープン回数を知らせるカウンタ(ディスク上ではなく、メモリ内のiノードテーブルにのみ存在)と混同してはならない。

5.6.5 ブロックキャッシュ

MINIXでは、システム性能を向上させるために、ブロックキャッシュを用いている。キャッシュはバッファの配列として実装される。各バッファはポインタを持つヘッダ、カウンタ、フラグ、そして1つのディスクブロック分の空間を持つ本体部分から構成されている。全ブロックは図5.33の様に最も最近使用したもの(MRU)から、最も昔に使用したもの(LRU)の順に、二重に連結されたリストを使ってつながれている。

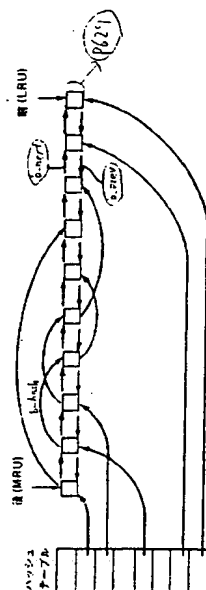


図 5.33 ブロックキャッシュの連結されたリスト

さらに特定のブロックがキャッシュ上にあるかを即座に判定する方法として、ハッシュテーブルを使用することもできる。ハッシュコードを持つすべてのブロックは、ハッシュテーブル内のエントリが指す単一に連結されたリスト上につながれる。現在のところ、ハッシュ機能は単にブロック番号から下位nビットを抽出しており、このため異なるデバイス上のブロックが同じハッシュチェーン内に存在している。

ファイルシステムがブロックを必要とする時は、手続き `get_block` を呼び出し、そのブロックのハッシュコードを算出して、ハッシュリストを検索する。ブロックが抽出された場合、ブロックのヘッダ内のカウンタが減少され、ブロックが使用中であることを示す。そのあとそれを指すポインタが返される。ブロックが抽出できなかった場合、LRUリストを検索し、キャッシュから除去するブロックを決定する。リストの前方のブロック(最も昔に使用されたブロック)がカウン

と上になっている場合、そのブロックが選択される。そうでなければ次のブロックが調べられ、といった具合である。ビットマップなど、ブロックによっては使用頻度に関わらず使用中に除去してはならないものもあるため、カウンタの検査は重要である。

除去するブロックを選択した後、そのヘッダ内の別のフラグが調べられ、ブロックを読み取った後で修正が行われたかどうかを確認する。修正が行われていればディスクに書き込まなくてはならない。そして必要なブロックを、ディスクタスクにメッセージを送信することによって読み取る。ファイルシステムはブロックの読み取りが終了するまで保留状態となり、終了した時点で実行を継続し、ブロックを指すポインタが呼出し例に返される。

ブロックを要求した手続きが処理を終了すると、別な手続き `put_block` を呼び出し、ブロックを解放する。`put_block` に対する引数のうちいくつかは、解放されるブロックのクラス(例えばiノード、ディレクトリ、データなど)を示している。クラスによって以下の2つの重要な決定がなされる。

- ブロックをLRUリストの前に置くか、後ろに置くか
- ブロックを即座に書き込むか否か(ブロックに修正が行われている場合)

二重連結ブロックなど、しばらく必要のないブロックはリストの前方に置かれ、次回空きバッファが必要になった時に要求できるようにしておく。すぐに必要になるブロックは、LRU形式に従いリストの最後に加えられる。

ディレクトリが修正された場合、クラッシュ発生時のファイルシステム破損の危険性を減らすため、即座にディスクへ書き込まれる。通常の修正済みデータブロックは、(1)LRUチェーンの先頭にきて除去されるか、(2)システムコールSYNCが実行された時、のうち、いずれかのイベントが発生するまで書き換えられない。

ブロックが修正されたことを示すヘッダ上のフラグは、ブロックの要求と使用を行ったファイルシステム内部の手続きによって設定される。手続き `get_block` と `put_block` は、連結されたリストの操作のみを行う。ファイルシステムがどのブロックをどの様な理由で必要としているかは、全く関知していない。

5.6.6 ディレクトリとパス

ファイルシステム内のもう1つ重要なサブシステムとして、ディレクトリとパス名の管理がある。OPEN など大半のシステムコールは、ファイル名を引数として持っている。必要であるのは、そのファイルのiノードであるため、ファイルシステムはディレクトリ・ツリーからそのファイルを検出し、iノードの位置を調べなくてはならない。

MINIXのディレクトリは1つのファイルに対して16バイトのエントリから構成されている。最初の2バイトは16ビットのiノード番号に、残りの14バイトはファイル名に使用されている。パス/user/ast/mbboxを検索する場合、システムはまずルートディレクトリ内からuserを検索

5章 ファイルシステム

し、次に /user 内の ast を検索し、最後に /user/ast 内の mbox を検索する。実際の検索はパスの要素それぞれに対し、1回に1つずつ行われる。図 5.11 はこれを示したものである。

これをさらに詳しく説明すると、標準的な MINIX の構成では /usr をフロッピーディスク (システムファイル) に、/user をフロッピーディスク 1 (ユーザーファイル) に用いている。以下の例では典型的なユーザーディレクトリとして /user/ast を使用している。

唯一複雑であるのは、マウントされたファイルシステムに遭遇した場合である。その仕組みを知るために、マウント方法を見てみる必要がある。ユーザーが以下の入力を入力を要求から行ったとしよう。

```
/etc/mount /dev/fd1 /user
```

フロッピーディスク 1 に格納されているファイルシステムは、ルート・ファイルシステムの /user にマウントされる。マウント前とマウント後のファイルシステムを図 5.34 に示す。

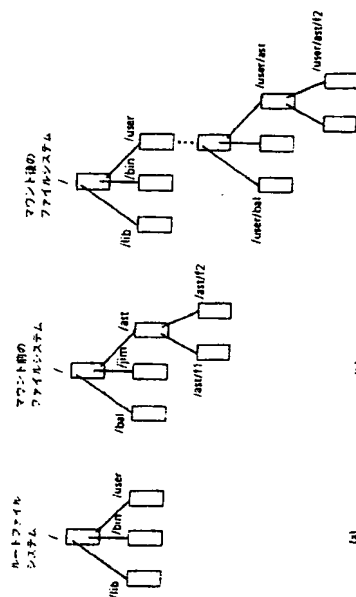


図 5.34 ファイルシステムのマウント
(a) ルート・ファイルシステム (b) マウント前のファイルシステム
(c) (b) のファイルシステムを /user にマウントした結果

マウント作業全体の鍵を握るのが、マウント成功後に /user の i ノードに設定されるフラグである。このフラグは i ノードにマウントが行われたことを示すものである。MOUNT システムコールは、新しくマウントされたファイルシステムのスーパーブロックを、super_block テーブルにロードし、そこに 2 つのポインタを設定する。さらに新しいファイルシステムのルートノードを i ノードテーブルに置く。

図 5.31 を見ると、メモリ内のスーパーブロックが、マウントされたファイルシステムに

5.6 MINIX ファイルシステムの概要

ある 2 つのフィールドを含んでいることがわかる。「マウントされたファイルシステムの i ノード」は、新しくマウントされたファイルシステムのルート i ノードを指すように設定される。次のフィールドである「マウントされた i ノード」は、マウントされた i ノード、つまりここでは /user の i ノードを指すように設定される。これら 2 つのポインタはマウントされたファイルシステムをルートに接続するために、その 2 つの間で「後着制」の働きをするものである (図 5.34 (c) では点線で示されている)。

/user/ast/12 というパス名を検索する場合、ファイルシステムは /user 内の i ノードのフラグを調べ、/user 上にマウントされたファイルシステムのルート i ノードを引き続き検索しなければならぬことがわかる。それではこのルート i ノードはどの様にして検出するのだろうか。

答えは簡単である。システムはメモリ内のスーパーブロックを、「マウントされた i ノード」フィールドが /user を指すものが検出できるまで検索する。これが /user にマウントされたファイルシステムのスーパーブロックである。スーパーブロックを検出すると、他のポインタに従ってマウントされたファイルシステムのルート i ノードを簡単に検出することができる。ここでファイルシステムを検索することができる。例えばフロッピーディスク 1 のルートディレクトリ内の ast を検索できる。

5.6.7 ファイル記述子

ファイルがいったんオープンされると、ファイル記述子がユーザープロセスに返され、以降の READ および WRITE コールで利用できるようになる。この項では、ファイルシステム内のファイル記述子の管理方法について考察していく。

カーネルやメモリマネージャ同様に、ファイルシステムはプロセステーブルの一部分をアドレス空間内に保存する。このうち、3 つのフィールドが特に重要である。最初の 2 つは 2 次元の 2 上りと 2 次元の 2 上りの i ノードを指すポインタである。図 5.11 に示したパス名の検索は、パス名が絶対的なものであるか、相対的なものであるかにより、ルートディレクトリかまたはワーキングディレクトリのどちらから検索を開始する。これらのポインタは、CHROOT および CHDIR システムコールを用いて、新しいルート、またはワーキングディレクトリを指すよう変更することができる。

プロセステーブル内で次に重要なフィールドは、ファイル記述子の番号をインデックスとした、配列である。これはファイル記述子が与えられた時に、そのファイルの場所を示すためのものである。単純に考えると、ファイル記述子 k に対するファイルを示すのみであれば、この配列の k 番目にそのファイルの i ノードを示す情報を持つだけで十分のように思われる。i ノードはファイルがオープンされるとメモリに取り出され、クローズされるまでそこに保存されるため、この方法が利用できる。

残念ながらこの様な単純な方法は、MINIX では (UNIX においても)、ファイルの所有者が行われる場合もあるためうまくいかない。次に読み書きを行うバイトを示すための 32 ビット長の値が、

それぞれのファイルに対してに用いられるため、問題が生じる。この値はファイル位置 (file position) またはファイルポインタと呼ばれ、LSEEK システムコールを用いて変更することができる。つまり問題となるのはファイル位置の保存場所である。

解決方法として考えられるのは、iノードにファイル位置を保存することである。残念ながら複数のプロセスが同時に同じファイルを開くと、それぞれのファイル位置を持つことになる。あるプロセスがLSEEKを実行し、次に別のプロセスが読み出しを行うような場合に問題が発生してしまう。したがって、iノードにはファイル位置を保存できないのである。

それではプロセステーブルに置けばどうか。ファイル記述子の配列と平行にもう1つ配列を設け、各ファイルの現在の位置を格納したらどうか。この方法も役に立たないのだが、その理由はいまひとつ明白である。基本的には、FORK システムコールの意味論が問題の原因となる。プロセスがフォークすると、親子のプロセスが各オープン・ファイルの現在位置を示す1つのポインタを共有する必要がある。

問題をわかりやすくするために、出力がファイルにリダイレクトされたシェルスクリプトの例を考えてみよう。シェルスクリプトが最初のプログラムをフォークすると、標準出力に対するそのファイル位置は0となる。この位置は子プロセスにも受け継がれ、例えば1Kの出力を書き込んだとすると、子プロセスが終了すると、共有するファイル位置は1Kになる。

ここでシェルスクリプトはさらにシェルスクリプトを読み取り、別な子プロセスをフォークする。2番目の子プロセスはシェルスクリプトからファイル位置として1Kを受け継ぐがなくてはならないため、最初のプログラムが終了した場所から書き込みを開始する。シェルスクリプトがその子プロセスとファイル位置を共有しなければ、2番目のプログラムは出力内容を最初のプログラムの追加するのではなく、上書きしてしまうことになる。

結論として、ファイル位置をプロセステーブルに置くことは不可能であることがわかった。MINIX では解決策として新しい共有テーブル flip を設け、すべてのファイル位置を保存することにした。その用法を図5.35に示している。ファイル位置を実際に共有させることにより、FORKは正しく実現され、シェルスクリプトも正しく動くことになる。

flip テーブルに格納しなければならないものは、共有のファイル位置であるが、iノードへのポインタも格納しておく必要がある。ここにiノードへのポインタが格納されていると、プロセステーブル内のファイル記述子配列内には、flip エントリを指すポインタだけを保存するのみでよい。flip エントリにはさらにそれを使用しているプロセスの数も含まれており、ファイルシステムがそのエントリを使用する最後のプロセスが終了した場合、そのスロットを再利用できるようにしている。

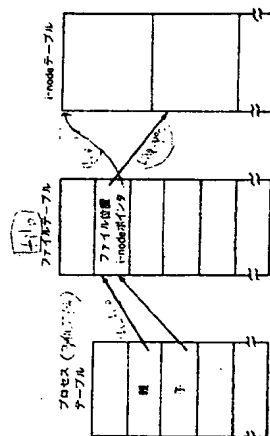


図5.35 親子間のファイル共有

5.6.8 バイブと特殊ファイル

バイブと特殊ファイルは、基本的な点で通常のファイルと異なっている。プロセスがディスクファイルから読み書きを行う時、操作は必ず2,300ミリ秒間のうちに終了する。最悪の場合いくつかのディスクアクセスを必要とするかもしれない。バイブから読み出す場合はこれとは異なっている。バイブが空であれば、読み出しを行おうとする者は、誰かがバイブ内にデータを送り込むまで待機しなくてはならない。場合によってはそれまで長時間も要することもある。同様に、読み出しから読み取りの場合も、プロセスは誰かが何かを入力するまで待たなければならない。

結論として、ある要求が終了するまで次の処理は行わないというファイルシステムの通常の規定は効力を持たないことになる。これらの要求をいったん中断し、後で再開することが必要である。プロセスがバイブから読み書きをする場合、ファイルシステムはバイブの状態を即座に調べ、操作が完了できるかを確認する。完了可能な操作であれば処理を行い、完了不可能であれば、ファイルシステムはシステムコールの引数をプロセステーブル内に保存し、時期を見てプロセスを再開できるようにする。

ファイルシステムは、特別な操作を行わずに呼び出し側をいったん停止させることができる。返答を行わなければ、呼び出し側は返答待ち状態のままブロックしてしまいうからである。したがってプロセスをいったん停止した後、ファイルシステムはメインループに戻り、次のシステムコールを待つ。別なプロセスがバイブの状態を変更し、保留中のプロセスを完了できるようにになった時、ファイルシステムはフラグを設定し、次のメインループで保留中のプロセスの引数をプロセステーブルから取り出し、コールを実行する。

増設や他のキャラクター特殊ファイルでは、若干状況が異なる。特殊ファイルのiノードには、2つの番号が格納されている。主デバイス番号と、副デバイス番号である。主デバイス番号は、デバイスクラスを示すものである (例えばRAMディスク、フロッピーディスク、ハードディスク、端末など)。それをファイルシステムテーブルに付するインデックスとして使用し、該当するタスク (すなわち入出力ドライバ) 番号に対応付ける。この様に、主デバイス番号によって呼び出し入

5.6 ファイルシステム

カドライブは決定するのである。副デバイス番号は、ドライブに対し引数として渡される。これは使用するデバイスとして、例えば端末2やドライブ1などを指定するものである。

プロセスが特殊ファイルからの読取りを行う時、ファイルシステムはファイルのiノードから主および副デバイス番号を取り出し、主デバイス番号を該当するタスク番号への対応付けのために、ファイルシステムテンプレートのインデックスとして利用する。タスク番号を得ると、ファイルシステムはタスクにメッセージを送る。メッセージには引数として副デバイス番号、実行する操作、呼出し側のプロセス識別子、バッファ・アドレス、そして転送するバイト数が含まれている。形式は、POSITION フィールドが使用されていないという点をのぞき、3章の図3.15と同じである。

ドライブが即座に作業を実行できる場合(例えば入力行がすでに端末から入力されている場合など)、データをドライブの内部バッファからユーザー用バッファにコピーし、作業の終了を伝える返答メッセージを送る。ドライブはデータをファイルシステムにはコピーしない点に注意されたい。ブロック型デバイスからのデータはブロックキャッシュを通過するが、キャラクタ型特殊ファイルはそこを通過しない。

ドライブが作業を実行できない場合、メッセージ・バッファ・メモータをその内部テーブルに記録し、即座にシステムコールを完了できないことを伝える返答メッセージをファイルシステムに送る。この時点で、ファイルシステムは遅延が空のバッファから読出しを行おうとしていることを発見した場合と同じ状況になる。プロセスが保留状態になったことを記録し、次のメッセージを待機する。ドライブが処理を完了できるだけのデータを得た時、それをブロック中のユーザーのバッファに転送し、ファイルシステムにその旨を伝えるメッセージを送る。ファイルシステムはここでユーザーに対し、ブロック解除のための返答メッセージを送り、転送されたバイト数を報告すればよい。

5.6.9 システムコール READ

後述するが、ファイルシステムの大半のコードはシステムコールを実行するために書かれている。したがって最も重要なシステムコールである READ の仕組みを述べることにし、以上の概要説明を補っていくことにする。

ユーザープログラムが以下の様なステートメントを実行し、通常のファイルを読み取ろうとした。

```
n=read(fd, buffer, nbytes);
```

ライブラリ関数 read は、3つの引数と共に呼び出されている。この結果、3つの引数を持つメッセージが、メッセージ型として READ に対するコードと共に作成され、メッセージがファイルシステムに送られる。返答を待ってブロックする。メッセージが到着すると、ファイルシステムはメッセージ型をユーザーに対するインデックスとして用い、読出しを行う手続きを呼び出す。

5.7 MINIX ファイルシステムの實現

この手続きはファイル記述子をメッセージから取り出し、それを `inproc` とし、読み出すファイルのiノード輸出に用いる(図5.35参照)。そして要求は1ブロック内に収まるように分割される。例えば現在のファイル位置が600で、1Kバイトが要求されている場合、要求は600から1023と、1024から1623の2つに分割される(1Kブロックと仮定して)。

これらを1つずつ格納し、そのブロックがキャッシュ内に存在するか調べる。存在しなければ、ファイルシステムは最も昔に使用したバッファのうちの現在使用されていないものを選択し、その内容が変更されていなければ書き直すように指示したメッセージをディスクタスクに送る。そして次にディスクタスクに読み取るブロックを要求する。

ブロックがキャッシュに入ると、ファイルシステムはデータをユーザー・バッファ内の該当する場所にコピーするよう指示したメッセージをシステムタスクに送る(すなわちバイト600から1023をバッファの先頭に、バイト1024から1623をバッファ内の424の位置に転送する)。コピーが終了したら、ファイルシステムは返答メッセージをユーザーに送り、コピーしたバイト数を伝える。

ユーザーがこのメッセージを受け取ると、ライブラリ関数 read は返答コードを取り出し、呼出し側に関数値としてそれを返す。

5.7 MINIX ファイルシステムの實現

MINIX ファイルシステムは比較的大きいが(Cで100ページを越える)、かなりわかりやすいものとなっている。システムコール実行のための要求が行われ、実行され、返答が行われる。本節では1ファイルごとにそれを見ていくことにしよう。特に重要と思われるところは解説を行っている。コード自体にも多くの注釈が含まれているので参照されたい。

5.7.1 ヘッドファイル

カーネルとメモリアネージャの間に、ファイルシステムにも各種のデータ構造とテーブルを定義するヘッドファイルが存在する。それではそのファイルシステムから見ていくことにしよう。ファイル `const.h`(7500行目)は、テーブルサイズやフラグなど、ファイルシステム全体で使われる定数を定義するものである。なかには `NR_BUFS` および `NR_BUF_HASH` の様に、システム性能をチューニングするために変更できるものもある。その他 `BOOT_BLOCK` および `SUPER_BLOCK` などは、性能とは関係がない。

次のファイル `buf.h`(7550行目)では、ブロックキャッシュの定義を行っている。配列 `buf` には、すべてのバッファが保存される。各バッファには `blockb`、そして多くのポインタ、フラグ、カウンタが含まれている。データ部分は5つの型のユニオンとして宣言されている(7565行目)。これは場合によってはブロックを文字列として、またはディレクトリなどとして参照するために効

果的にからである。

バックアップ3のデータ部分を文字配列として参照する場合、buf[3]は共用体全体を参照し、そこからファイルb_dataを選択するためには、buf[3].b_dataとするのが正しい方法である。この構文は正しいが、多少面倒であるため、7588行目ではマクロb_dataを定義し、buf[3].b_dataと略記する。b_data(ユニオンのフィールド)では2つの下線を適用しているが、b_data(マクロ)では略記しやすいように1つしか使っていない点に注意されたい。ブロックの別な部分をアクセスするためのマクロは、7588から7592行目に含まれている。

このファイルでいう1つ注目したいのは、すべての配列と変数に対してEXTERNを用いていることである。このファイルがコードファイルに含まれている時、EXTERNは0009行目で定義されているように、externに置き換えられる。しかしファイルtablecでは、保存領域の割当てを行なうためにヌルストリングとして定義されている。Cでは(Kernighan and Ritchie, 1978, 頁p.206)、大域変数は1つのファイルを除くすべてのファイルでexternとして宣言されなくてはならないことになっている。この点を理解していないコンパイラやプログラマが多いのは残念なことである。これと同じような問題がカーネルとメモリマネージャにも見られた。

ファイル最後のマクロ(7601から7610行目)は、ブロックの種類を定義するものである。ブロックが使用後にバックアップ・キャッシュに戻された時、これらの値のいずれかが与えられ、キャッシュマネージャはそれによってそのブロックをLRUリストの先頭または後ろのどちらに置くか、またディスクにすぐ書き込みを判断する。

ファイル(7611から7650行目)は、dmapマッピングの定義を行うものである。マッピング自身は初期値と共にtablecで宣言されており、複数のファイルに含まれないようになっている。devhが必要なる理由はここにある。マッピングは主デバイス番号と該当するタスクの対応付けを行う。

ファイル(7651から7700行目)は現在のファイル位置と、iノードのポインタを格納するための中間マッピングを持つ(図5.35)。さらにファイルが読み取り、書き込み、あるいはその両方、のいずれを行なうためにオープンされたのか、そして現在そのエントリを指しているファイル記述子の数を格納する。

プロセスマッピングのファイルシステム部分は(7750行目)に記述されている。モードマッピングは現在のルートディレクトリとマッピングエントリに対するiノードを指すポインタ、ファイル記述子の配列、uid、gid、および端末番号が含まれている。残りのフィールドは例えは空のバリエーションを認めるようにして途中で保留になっていたシステムコールの引数を保存するために使われる。フィールドfp_suspendedおよびfp_revivedは、実際には1ビットしか必要としないが、大半のコンパイラはビットフィールドよりも文字に対して質の良いコードを作成するため、この様にしている。

次は大域変数を持つファイル(7800行目)である。受け取るものと返答メッセージに対するメッセージ・バックアップも、システムスタックと共にここに存在する。MINIXのブート後、ファイルシステムが初めて起動されると、非常に小さなアセンブラ手続がスタックポインタを(istack)の最上部

に設定する。

次はinodeヘッダフィールドである(7850行目)。何度も触れたように、ファイルをオープンすると、そのiノードがメモリ内に読み込まれ、ファイルをクローズするまでそこに保存される。これらのiノードはこのフィールドに保存される。以上で、大半のフィールドの意味は理解できたはずである。ただしiseekに関しては、若干の説明を加えておこう。最適化の一環としてファイルシステムは、ファイルが連続して読み出されることに気づくと、要求がくる前にブロックをキャッシュに読み込もうとする。ランダムアクセス・ファイルでは、先読みの機能は提供されていない。LSEEK コールが実行されると、フィールドiseekが設定され、先読みを禁止する。

フィールドparam(7900行目)は、メモリマネージャ内の同じ名前のフィールドと似ている。引数を含むメッセージ・フィールドの名前を定義し、例えばメッセージ・バッファの1フィールドを選択するmmlplの代わりに、bufferとして参照できるようにする。

super(7950行目)ではスーパーブロック・テーブルの宣言が行われている。システムがブートされると、ルートデバイスのスーパーブロックがここへロードされる。ファイルシステムのマウント時にも、そのスーパーブロックがここに配置される。

最後に型の定義が(8000行目)で行われている。2つの型、すなわちディレクトリ・エントリとディスタント・ノードが定義されている。

5.7.2 テーブル管理

主要な各テーブル、すなわちブロック・iノード、スーパーブロックなどのテーブルを管理する手続は、それぞれが1つのファイルにまとめられている。これらの手続はファイルシステムその他の部分でも頻繁に使用され、テーブルとファイルシステム間の主要なインターフェイスとなっている。したがって、そのファイルシステムコードから学習していくことにしよう。

■ ブロック管理

ブロックキャッシュはファイルcacheの手続によって管理されている。このファイルには図5.36で示した5つの手続が含まれている。最初の手続はget_block(8079行目)はファイルシステムによるデータブロック取得方法として標準的なものである。ファイルシステム手続がユーザーデータブロック、ディレクトリ・ブロック、スーパーブロック、その他の各種のブロック

get_block	読み書きブロックを得る
put_block	読みのget_blockによって取得されたブロックを返す
alloc_block	新しいiノードを割り当てる (ファイル成長)
free_block	ゾーンの解放 (ファイルの削除)
rev_block	ディスク・キャッシュ間のブロック転送
invalidate	デバイス間の空をキャッシュを一時停止

図 5.36 ブロック管理の手続

5.7 MINIX ファイルシステムの实现

近いゾーンを検索し、ファイルのゾーンをひとまとめにしようとすると、ビットマップ内のビット番号と、ゾーン番号の対応は 8268 行目で行われている。ビット 1 が最初のデータゾーンに相当することがわかる。

ファイルを削除すると、そのゾーンがビットマップに返却される。[free_zone] (8275 行目) はこれからゾーンの返却作業を行う。free_bit を呼び出し、free_zone とビット番号を引数として渡すだけである。free_bit は、空の i ノードを返却する際にも用いられるが、もちろんその場合は i ノードマップを最初の引数として用いる。

キャッシュの管理には、ブロックの読み書きが必要である。ディスクとの単純なインターフェイスのために、手続は [rw_block] (8295 行目) が用意されている。この手続は単一ブロックの読み書きを行うものである。同時に手続は rw_inode と rw_super が、それぞれ i ノードとスーパーブロックの読み書き用に用意されている。

このファイルの最後の手続は [invalidate] (8326 行目) である。これはディスクがマウント解除される時に使用され、マウント解除されたファイルシステムに属するすべてのブロックをキャッシュから除去する。この作業を行っておかないと、デバイスが再度利用された時(別のフロッピーディスクによって)、ファイルシステムは新しいディスクのブロックではなく、古いものを使ってしまいう可能性がある。

■ i ノード管理

ブロックキャッシュ以外に手続の助けを必要とするテーブルがある。i ノードテーブルがそれである。多くの手続は、ブロック管理の手続きと機能的に同じである。図 5.37 に一覧を示す。

get_inode	i ノード内の inode を得る
put_inode	不要な inode を返す
alloc_inode	新しい inode を割り当てる(新しいファイル)
write_inode	inode 内のいくつかのフィールドのクリアする
free_inode	inode を解放(ファイルの削除)
rw_inode	i ノードとディスク間の inode を伝送する
flag_inode	i ノードが使用中であることを示す

図 5.37 i ノード管理用手続き

手続は [get_inode] (8379 行目) は get_block に似ている。ファイルシステムの任意の部分が i ノードを必要とした場合、get_inode を呼び出し、それを得る。get_inode はまず i ノードテーブルを検索し、i ノードの存在を確認する。存在が確認できたら、使用カウンタを増加してそれを指すポインタを返す。この検索作業は 8389 から 8406 行目に記述されている。i ノードがメモリ内に存在しない場合は、rw_inode を呼び出し、i ノードをロードする。

i ノードを必要とする手続が終了した後、手続は [put_inode] (8421 行目) を呼び出して i ノード

5 章 ファイルシステム

を読み取る場合、デバイスとブロック番号を指定して get_block を呼び出す。

get_block を呼び出すと、最初にブロックキャッシュを調べて要求したブロックが存在しているかどうかを確かめる。もし存在していれば、それを指すポインタを返す。存在しない場合には、ブロックを読み出さなければならない。キャッシュ内のブロックは連結されたリスト NR_BUF_HASH (32) 上につながれている。それぞれのリストにつながれているブロックはすべてブロック番号の最後が同じ 5 ビットの並び、すなわち 00000, 00001...11111 を持っている。

8099 行目のステートメントは、要求したブロックがキャッシュ内に存在するならば、そのブロックが属していることになっているリストの先頭を bp が指すよう設定する。8101 行目のルーブはこのリストを検索し、ブロックを検出する。存在すれば 8106 行目で呼び出し側に返される。

ブロックがリスト上に存在しない場合は、キャッシュ内にも存在しないことになるため、最後に使用されたキャッシュのうち、現在使用されていないものが用意される。ビットマップなど、現在使用されているブロックを除去してはならない。選択したバッファはハッシュチェーンから削除される。新しいブロック番号で使用するために別のハッシュチェーン上に移すためである。内容が修正されている場合は 8139 行目でディスクに書き込まれる。

バッファが使用可能になると、即座に新しい引数が設定され、ブロックがディスクから読み取られる。これには例外が一つある。ファイルシステムが一つのブロック全体を書き込もうとする場合、最初に古い内容を読み出すのは無駄である。この場合は、ディスクの読み取りが省略される(8149 行目)。新しいブロックが読み取られた後、get_block は呼び出し側にそれを指すポインタを返す。

ファイルシステムがファイル名の検索のために、一時的にディレクトリ・ブロックを必要としていてと仮定しよう。get_block を呼び出してディレクトリ・ブロックを得る。そのファイル名を検出した後、put_block を呼び出して(8157 行目)ブロックをキャッシュに返却する。これにより、別なブロックにおいて必要になった時、このバッファを利用できることになる。

手続は put_block は、新たに返却されたブロックを LRU リスト上に置き、場合によってはディスクに書き込みを行う。まず最初に(8189 行目)ブロックを LRU リストの現在位置から除去する。次にブロックの種類を要するために呼び出し側が指定したフラグ block_type を基に、LRU リストの先頭、または最後にそれを置く。しばらく必要ないと思われるブロックは先頭に置かれ、そこで再び利用される。すぐに必要となるブロックは最後に置かれ、しばらくそこで待機する。

ブロックが LRU リスト上に再配置されると、ディスクにブロックの書き込みをすぐに行う必要があるかが検査される(8057 から 8225 行目)。i ノード、ディレクトリ・ブロック、その他ファイルシステムの機能自体に関連のあるブロックはこれに属するため、その時点でディスクに書き込まれる。

ファイルが成長すると、場合によっては新しいゾーンを割り当て、新しいデータの保存に備えてくなくてはならないことがある。手続は [alloc_zone] (8235 行目) は、新しいゾーンの割当てを行うものである。ゾーン・ビットマップを検索し、空きゾーンを検出する。現在のファイルのゾーン 0 に

を返す。これにより、使用カウンタ `i_count` を減少することになる。ここでカウンタが0になった場合、ファイルが使用されなくなることになり、iノードはテーブルから削除される。内容が変更されている場合はディスクに書き込まれる。

`i_link` フィールドが0の場合は、ファイルを指しているディレクトリ・エントリが存在しないことになり、そのゾーンがすべて解放される。使用カウンタが0になると、リンク数が0になるとはその意味も、原因も、結果も全く異なる点には注意されたい。

新しいファイルが作成される場合、iノードを割り当てなくてはならない。この作業は `alloc_inode(8446 行目)` によって行われる。ゾーンの解放はファイルごとに近隣のゾーンを割り当てようとしたが、これと異なり任意のiノードを使えばよい。

iノードが得られたら、`get_inode` を呼び出して、iノードをメモリ内のテーブルに取り出す。次にそのフィールドの一部が8486 行目で、また残りの部分は `wipe_inode(8503 行目)` で初期化される。ファイルシステムの他の部分で、特定のiノード・フィールド(全部ではない)をクリアするために `wipe_inode` が必要となるので、この様な作業分組が行われている。

ファイルを削除すると、`free_inode(8525 行目)` が呼び出され、そのiノードが解放される。ここではiノード・ビットマップ内の該当するビットが1に設定されるだけである。手続き `rw_inode(8543 行目)` は、`rw_block` と同じように、ディスクからiノードを取り出す役を担っている。以下の手順で作業を実行する。

- ① 必要なiノードを含むブロックを算出する。
 - ② `get_block` を呼び出してブロックを読み取る。
 - ③ iノードを抽出して、それを `inode` テーブルにコピーする。
 - ④ `put_block` を呼び出して、ブロックを返却する。
- 手続き `dup_inode(8579 行目)` は、iノードの使用カウンタを増加するだけである。

■ スーパーブロックの管理

ファイル `super.c` には、スーパーブロックとビットマップの管理を行う手続きが記述されている。このファイルには、図 5.38 で示している7つの手続きがある。

<code>load_bit_maps</code>	別のファイルシステムのビットマップを得る
<code>unload_bit_maps</code>	ファイルシステムのマウント解除時のビットマップを返す
<code>alloc_bit</code>	ゾーン・i-node マップからのビットを割り当てる
<code>free_bit</code>	ゾーンおよびi-node マップ内のビットを解放する
<code>get_inode</code>	ファイルのためのスーパーブロック・テーブルを返す
<code>set_inode</code>	ゾーンとブロックを管理するに用いられるビット数
<code>scale_factor</code>	メモリとディスク間のスーパーブロック転送

図 5.38 スーパーブロックとビットマップ管理用の手続き

`load_bit_maps(8631 行目)` はルートデバイスがロードされる時、または新しいファイルシステムがマウントされる時に呼び出される。すべてのビットマップ・ブロックを読み取り、スーパーブロックがそれらを指すように設定する。スーパーブロック内の配列 `s_imap` と `s_zmap` は、それぞれiノード・ビットマップ・ブロックとゾーン・ビットマップ・ブロックを指すものである。

ファイルシステムがマウント解除されると、そのビットマップが `unload_bit_maps(8669 行目)` によってディスクに書き込まれる。

iノードまたはゾーンが必要な場合は、前述の様に `alloc_inode` または `alloc_zone` を呼び出す。これらは両方とも `alloc_bit(8689 行目)` を呼び出し、実際に関連するビットマップを検索する。検索には以下の三重のループが使用される。

- ビットマップの全ブロックに対する外部ループ
- ブロックの全ワードに対する中間ループ
- ワードの全ビットに対する内部ループ

中間ループは現在のワードが0の補数、すべてのビットが1であるワードに等しいことを調べ、等しければ、空iノードもゾーンも存在しないことになるため、次のワードが調べられる。異なる値を持つワードが検出されたら、少なくとも1つの0のビットが含まれていることになるので、内部ループに入り、空(すなわち0)ビットを検出する。すべてのブロックを調べたが該当するものがなかった場合は、空のiノードもゾーンも存在しないことになり、コード `NO_BIT(0)` が返される。

ビットの解放には検索が不要なため、割り当てより簡単である。`free_bit(8747 行目)` は、解放するビットを持つビットマップ・ブロックの算出を行い、そのビットを0に設定する。ブロック自体は常にメモリ内にあり、スーパーブロック内の `s_imap` または `s_zmap` ポインタによって指し示されている。

次の手続き `get_super(8771 行目)` は、特定のデバイスに対するスーパーブロックを検出する。例えばファイルシステムのマウントを行う場合、そのデバイスはまだマウントされていないことを確認しなければならない。この検査は `get_super` によりデバイスを検出することで実現する。デバイスが検出できなかった場合は、そのファイルシステムがまだマウントされていないことになる。ブロックとゾーンの要領はブロック番号を左に、そしてゾーン番号を右にシフトして行う。シフトする量はゾーンあたりのブロック数に依存する。この数はファイルシステムによって異なる。手続き `scale_factor(8810 行目)` はこの値を求める。

最後に `rw_super(8824 行目)` は、前述した `rw_block` と `rw_inode` に似ており、スーパーブロックの読み書きを行うために呼び出される。

■ ファイル記述子の管理

MINIX ではファイル記述子と `filep` テーブルを管理するために、特別な手続きを記述している(図 5.35 参照)。これらは `filep` の記述に記述されている。ファイルが作成、またはオープンされると、そのファイル記述子と空の `filep` スロット上が必要になる。手続き `get_filep(8871行目)` を使ってこれらを抽出することができる。CREATE や OPEN が成功すると確認できるまでは多くの検査が必要のため、使用中の印は付けられていないので注意されたい。

`get_filep(8816行目)` はファイル記述子が指定された範囲内にあるかを検証し、その `filep` ポインタを返す。

このファイルの最後の手続きは `find_filep(8930行目)` である。プロセスによる破損したパイプ(すなわち他のプロセスが読取り用にオープンしていないパイプ)への書き込みを検知するために必要である。filep テーブルを検査して、読取りを行う可能性のあるものを検出する。

5.7.3 メインプログラム

ファイルシステムのメインループは、8992行目から始まっているファイル `main.c` に記述されている。その構造は、メモリアネージャと入出力アスクのメインループに非常によく似ている。`(get_work)` を呼び出すと、次のメッセージの到着を待つことになる(パイプまたは端子上でそれまで保留状態にあったプロセスが処理される場合を除く)。大域変数 `who` を呼出し、そのプロセス番号と `filep` スロット番号に設定し、もう1つの大域変数 `is_call` に、実行するシステムコール番号を設定する。

メインループに戻った時、3つのフラグが設定される。fp は呼出し側のプロセス番号、super_user は呼出し側がスーパーユーザーか否かを示し、dont_reply は FALSE に初期化される。次に最も重要な作業が行われる。システムコールを実行する手続きを呼び出す。この手続きは `is_call` を、手続きへのポインタ `call_vector` 配列に対するインデックスとして使用し、その手続きを呼び出す。

メインループに制御が戻った時、`dont_reply` が設定されていなければ返答は行われない(例えばプロセスが他のパイプを読み出すようにしてブロックしてしまつた)。それ以外の場合は、なんらかの返答が行われる。メインループの最後のステートメントは、連続的に読み取られるファイルを検出するために設けられたものであり、実際に必要がある前に次のブロックをキャッシュに読み込み、システムの性能を向上させる。

手続き `get_work(9016行目)` は、これまでにブロックされた手続きのうち、再開できるものがないかを調べる。再開できる手続きがあれば、新しいメッセージより高い優先度を与えられる。特に内部で行う作業がない場合、ファイルシステムはメッセージを受け取るためにカーネルを呼び出す(9042行目)。

システムコールが正常終了しても、または異常終了しても、reply(ライン 9053行目)を用いて呼出し側に返答を行う。基本的には send が失敗することはないが、カーネルは確認のために状態

コードを返す。

ファイルシステムが実行を開始する前に、`filep` (9069行目)によって自身の初期化を行う。この手続きによってブロックキャッシュが用いられる連番されたリストが構築され、64Kの境界を超えないようにパイプをすべて削除する(IBM PC の DMA チップは 64Kの境界を超えることができないため)。そしてアポートディスクから RAM ディスクをロードし、スーパーブロックテーブルを初期化し、ルートデバイスのスーパーブロックと i ノードを読み取る。すべてに問題がなければ、i ノードとゾーン・ビットマップがロードされる。最後にそれらの数値が正当であることを確認する。

アポートディスクを作成する時、RAM ディスクのビット単位のコピーが MINIX パイナリ の後に含まれる。手続き `load_ram` は、最初に各種設定(RAM ディスクドライバに RAM ディスクの配置先と、その大きさを伝えることも含まれている)を行い、ブロック単位で RAM ディスクにコピーする。

■ ディスバッチャ・テーブル

ファイル `(alloc)(9300行目)` は、どの手続きがどのシステムコール番号を処理するかを決定するために、メインループで使用するポインタの配列を含んでいる。同様のテーブルがメモリアネージャ内にも存在している。

ただし 9416行目のテーブル `dispatch` は、新たに設けられたものである。このテーブルは 0 から始まる1つの列を各種のデバイスに対して持っている。デバイスをオープン、クローズ、読み書きすると、操作を行う手続きの名前がこのテーブルによって提供される。これらの手続きはすべてファイルシステムのアドレス空間に位置している。これらのうち、大半は特に何も行わないが、なかには実際に入出力を行うタスクを呼び出すものもある。各種デバイスに該当するタスク番号も、このテーブルが提供する。

新しく主デバイスが MINIX に追加されると、このテーブルに1行が追加され、デバイスのオープン、クローズ、読み書きの際の処理について指示する。簡単な例として、テープドライブが MINIX に追加されたとき、その特殊ファイルがオープンされた時に、テーブル内の手続きによりテープがすでに使用中であるかどうかを調べることができる。

5.7.4 ファイルに対する操作

この項では、ファイル(ディレクトリではなく)に操作を行うシステムコールについて考察する。ファイルの作成、オープン、クローズからその読み書きまでを調べていく。

■ ファイルの作成、オープン、クローズ

ファイル `open.c` は5つのシステムコールに対するコードを含んでいる。CREATE、MKNO、OPEN、CLOSE、LSEEK がそれである。これらをつつ見ていくことにしよう。ファイルの

作成は以下の手順で行われる。

- ①新しいファイル割当て、初期化する。
- ②対応するディレクトリに新しいファイルを追加する。
- ③新しいファイルのファイル記述子を設定し、戻却する。

CREATE を処理する手続きは `do_create(9479 行目)` である。メモリーマネージャの場合と同様に、システムコール XXX は手続き `do_XXX` によって実行されるという規則をファイルシステム内で用いている。

`do_create` は、最初は新しいファイルの名前を取り出し、^(5.5) 空のファイル記述子と `flip` テーブルスロットが使用可能であるかどうかを確認する。新しい `i` ノードを実際に作成するのは手続き `new_node` であり、これは 9496 行目で呼び出されている。`i` ノードを作成できなかった場合、`new_node` が大抵変数 `err_code` を設定する。

`do_create` が行う作業は、ファイルがすでに存在するか否かによって異なってくる。ファイルが存在しない場合、9504 から 9521 行目はスキップされ、テーブルスロットが要求されて、ファイル記述子が返される。

ファイルが存在する場合は、ファイルシステムはファイルの種類やモードなどを確認しなければならない。通常のファイルに対して CREATE を行くと、そのファイルは長さ 0 に切り捨てられる。書込み可能な特殊ファイルに対して CREATE を実行すると、書込みのためにオープンされることになる。それをディレクトリに対して実行すると、必ず拒否されてしまう。

`do_create` のコードや、他の多くのファイルシステム手続きでは、各種のエラーや不正な状態を検出するためのコードがたたくと記述されている。これらのコードはエラーのない、強固なファイルシステムを作るには欠かせないものである。何も問題がなければファイル記述子と先頭に配置された `flip` スロットが、ここで割当て済みとして印付けられ、ファイル記述子が返される。最初から割当て済みとは印付けられず、必要に応じて途中で終了しやすいうにしている。

MKNOD システムコールは、`do_mknod(9541 行目)` によって処理される。この手続きは `do_create` と似ているが、`i` ノードを作成し、そのためのディレクトリ・エントリを設けるだけである。`i` ノードがすでに存在する場合、システムコールはエラー終了する。`do_create` で見たような、状況に応じた分析は必要ない。

`i` ノードの割当てとファイルシステムへのパス名設定は、`new_node(9557 行目)` によって行われる。9575 行目のステートメントは、パス名の解析(すなわちパス名の要素ごとに調べる)を最後のディレクトリまで行う。3 行後の `advance` に対する呼び出しは、パス名の最後の要素をオープンできるかどうか調べるものである。

```
id=create("/user/ast/foober", 0755);
```

例えば上記の操作を呼出しでは、`last_dir` は `/user/ast` の `i` ノードをテーブル内にロードし、それ

に対するポインタを返す。ファイルが存在しない場合、この `i` ノードを使ってディレクトリに `foobar` を追加することになる。ファイルの追加、削除を行う他のシステムコールも、すべて最初にパス内の最後のディレクトリをオープンするために `last_dir` を使用する。

`new_node` によってファイルが存在しないことが判明した場合、9581 行目で `alloc_inode` を呼び出し、新しい `i` ノードの割当てとロードを行い、それを指すポインタを返す。空の `i` ノードが残っていない場合、`new_node` は失敗し、`NIL_INODE` が返される。

`i` ノードが割り当てられれば、9591 行目で作業を継続し、ファイルのいくつかを埋め、ディスクに再度書き込み、ファイル名をパス名中の最後のディレクトリに書き込む(9596 行目)。ここでもファイルシステムは飽えずエラーを調べ、検出時には使用している `i` ノードやブロックなど、すべての資源を安全に解放する。例えば `i` ノードが不足した時、現在のコールの結果を取り消したり、呼出し側に対してエラーコードを返さずに、あえて MINIX をパニック状態にするつもりならば、ファイルシステムはかなり簡便化できるだろう。

次の手続きは `do_open(9622 行目)` である。各種検査を行った後、`eat_path` を呼び出し、ファイル名の解析を行い、`i` ノードをメモリー内にフェッチする。`i` ノードが使用可能になると、モードを調べてファイルがオープン可能かどうかを確認する。9645 行目における `forbidden` の呼出しでは、`rxw` ビットを検査を行う。ディレクトリと特殊ファイルは後で処理される。最後にファイル記述子が関数値として返される。

ファイルのクローズは、ファイルのオープンよりも簡単である。`do_close` によってこの作業が行われる(9680 行目)。バイブと特殊ファイルは若干注意しなければならないが、通常のファイルでは `flip` カウンタを減少し、0 になった場合は `i` ノードが `put_inode` を使用して戻される。

`i` ノードを戻すことにより、`i` ノードテーブル内のカウンタは減少される。したがって結果的にテーブルから削除されることになる。この操作は `i` ノードの解放(`i` ノードが未使用であることをビットマップに設定すること)とは違う。`i` ノードはファイルがそれを持つすべてのディレクトリから除去された場合に限り、解放することができる。

このファイルの最後の手続きは `do_lseek(9721 行目)` である。シークが実行されると、この手続きが呼び出され、ファイル位置を新しい値に設定する。

■ ファイルの読出し

ファイルをオープンすると、それを読み書きすることができる。最初に読出しについて説明し、次に書込みについて説明する。いくつかの点でこの 2 つの作業は異なっているが、`do_read(9784 行目)` も `do_write(10125 行目)` はどちらも共通の手続き `read_write(9794 行目)` を呼び出し、そこで大半の作業を行うという点では一致している。

9811 から 9818 行目までのコードはメモリーマネージャが、ファイルシステムを使ってそのユーザー空間にセグメント全体をロードする際に用いる。通常のコールは 9821 行目から処理を開始し、そこでは正当性検査(例えば書込みだけが行えるファイルから読み取ろうとしたなど)と変数

5.3 ファイルシステム

の初期化が行われる。キャラクター特殊ファイルからの読出しはブロックキャッシュを通過しないため、9836行目で削除されている。

9844から9854行目までの検査は、デバイスの許容量よりも成長する可能性のあるファイルや、ファイル終端を超えて書込みを行うことによってファイル内に空間を作ってしまったものなど、書込みに関するものである。これはMINIXの概要説明で触れたように、1つのゾーンに複数のブロックが存在すると、問題を問題を引き起こすことになるからである。バリエーションも特殊であり、検査が必要である。

少なくとも通常のファイルにおいては、読出し機構の主な処理となるのが9861行目から始まるループである。このループは要求をいくつかの処理単位(chunk)に分割し、それぞれが1つのディスクブロック内に収まるようにする。処理単位は、現在位置から始まり以下のいずれかの状況になるまでである。

- すべてのバイトを読み出した場合①
- ブロック境界に達した場合②
- ファイル終端に達した場合③

この様な規定は、1つの処理単位がいくつかのディスクブロックにまたがらないことを意味している。図5.39では、処理単位のサイズを決定する方法として3つの例を紹介している。実際の処理は9862から9871行目で行われている。

処理単位の実際の読出し作業はrw_chunk(9874行目)によって行われる。制御が戻された後、各種カウンタとポインタが増加され、次の操作を開始する。ループが終了した後、ファイル位置と他の変数を更新することができる(例えばバイブ・ポインタなど)。

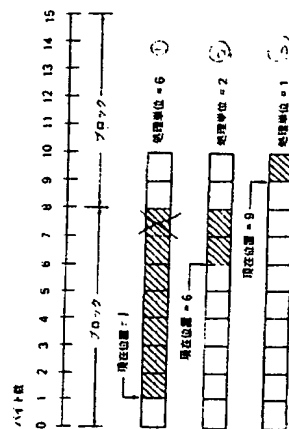


図5.39 10バイトのファイルに対して処理単位のサイズを決定する3つの例
ブロックサイズは8バイトであり、必要なバイト数は6である。処理単位は8
バイトである。

5.7 MINIX ファイルシステムの實現

先読みを指定した場合は、読み取るiノードと位置が次第変数に保存され、返答がユーザーに送られた後、ファイルシステムが次のブロックの処理を開始できるようにしておく。通常、ファイルシステムは次のディスクブロックを待ってブロックし、その間にユーザープロセスはすべてに与えられているデレークの処理を行うことができる。

手続rw_chunk(9919行目)は、iノードとファイル位置を指定し、それらを物理ディスクブロック番号に変換し、そのブロック(またはそのブロックの一部)をユーザー空間に移動するよう指示する。相対的なファイル位置と物理ディスクアドレスの対応付けは、iノードと間接ブロックを聞いているread_mapが行う。通常のファイルの場合は、9944と9945行目の変数およびdevに物理ブロック番号と、デバイス番号がそれぞれ格納されることになる。9964行目で呼び出しているget_blockは、キャッシュハンドラに対してブロックの検出を依頼するものであり、必要に応じてそれを読み取る。

ブロックに対するポインタを得ると、9972行目のrw_userの呼出しにより、要求した部分をユーザー空間に転送する。そしてブロックをput_blockによって解除し、後で必要に応じてキャッシュから削除できるようにする(get_blockで要求すると、ブロックのヘッド内にあるカウンタはそれが使用中であることを示し、削除の対象からははずす。そしてput_blockはカウンタ値を減少させる)。

read_map(9984行目)はiノードを調べることによって、論理ファイル位置を物理ブロック番号に変換する。ファイルの先頭から最初のiゾーン(iノード内に含まれている)のうちのいずれかに存在するブロックに関しては、どのゾーンが必要であり、そしてそれがどのブロックであるかの決定は、簡単な処理で行える。それ以上のブロックを使用しているファイルに関しては、1つ以上の間接ブロックを読み取らなければならない。

手続rw_user(10042行目)は、単にシステムタスク用にメッセージを作成し、送信するものである。実際のコピー作業は、カーネルが行う。ファイルシステムがコピーを行うことはない。ファイルシステムはユーザーがメモリ内のどこに位置しているのかを知らず知らない。この様な余分なオーバーヘッドは高度にモジュール化されたシステムで、私わなくてはならない犠牲である。

最後にread_ahead(10082行目)は論理位置を物理ブロック番号に変換し、get_blockを呼び出してブロックがキャッシュ内に存在することを確かめ、ブロックを即座に戻す。結局、ブロックに対して何も行わないのである。ブロックをすぐに利用する時のために、ブロックをできるだけキャッシュ内に置くというものである。

read_aheadは、main内のメインループでしか呼び出せない点に注意されたい。READシステムコールの処理の一部として呼び出すことはできない。read_aheadの呼出しは必ず返答を終えてから行い、先読み中にファイルシステムがディスクを待機しなければならない場合でも、ユーザーが実行を継続できるようにする。図5.40はファイルの読出しに関する主な手続きの関係を図示したものである。

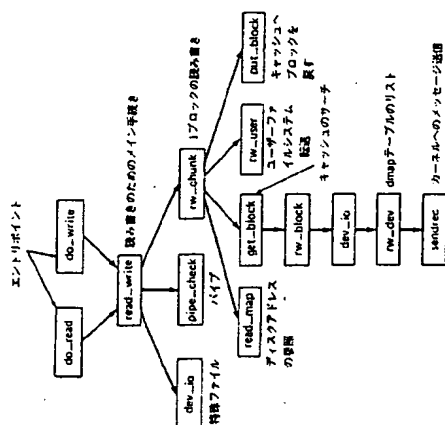


図5.40 ファイルの読出しに関するいくつかの手続き

■ ファイルの書き込み

ファイルの書込みは、読取りと似ている。ただし書込み作業では、新しいディスクアブロックを割り当てなければならぬ。`write_map[10135]`行目は `read_map` と似ているが、`i` ノードと関係するアブロック内の物理アブロック番号を検索する代わりに、新しいアブロック番号を `i` ノードに格納する。正確にいえば、アブロック番号ではなく、ゾーン番号を格納する。

`write_map` のコードはいくつかの状況に対応するため、若干長めで、細くなっている。ソースがファイルの先頭に近ければ、iノードに格納される(10160行目)。

一番面倒なのは、ファイルが単一間接ブロックで取扱うことのできる最大の大きさとなった場合である。この場合二重間接ブロックが割り当てられる。次に単一間接ブロックの割当てが必要でなくなり、そのアドレスが二重間接ブロックに置かれる時に、二重間接ブロックの割当ては成功した。単一間接ブロックを割り当てることができない場合(すなわちディスクが満杯の場合)、ビットマップを確認しないように、注意深く二重間接ブロックを解放しなければならぬ。

ここでも、バッキング状態となった構文なのであれば、コードはかなりの簡素化できる。しかしユーザーの立場からいえば、ディスク空間が足りなくなっても、コンピュータのクラッシュや、ファイルシステムのエラーの増殖を発生せずに、WRITEからのエラーを返すにはいかにだけ助かるかわからない。

writec 内の次の手続きは `clear_zone` である。これはファイル終端を超えてシークを行った場合に、データを書き込む前に実行され、ブロックの消去を処理する手続きである。幸運にも、こ

5.7 MINIX ファイルシステムの実現

の様な状況はそれほど頻繁には発生しない。

new_block (10265 行目)は、新しいブロックが必要になった場合には必ず `rw_chunk` 内に記述されている 9955 行目で呼ばれる、図 5.41 はシリアライズされたファイルの成長を 5 段階に分けて示したものである。ブロックサイズは 1K で、ゾーンサイズは 2K となっている。

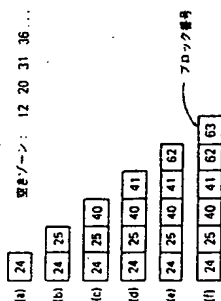


図5.41 シーケンシャル・ファイルの成長

(a)-(1) 1Kブロックの連続割当て。ゾーンは2Kである。

new_block が初めに呼び出された時は、ゾーン 12(ブロック 24 と 25)を割り当て、次に呼び出された時は、すでに割り当てられているがまだ使用されていないブロック 25 を使用する。3 回目の呼出しでは、ゾーン 20(ブロック 40 と 41)が割り当てられ、といった具合になる。zero_block (10318 行目)はブロックをクリアし、それまでの内容を削除する。

■ バイフ

パイプも通常のファイルと多くの共通点を持つ、しかしこの項では相違点に焦点を当てることとする。まずパイプは CREAT ではなく、PIPE システムコールによって作成される。PIPE コールはファイル `pipe.c` 内の `do_pipe` (行目 1184) によって処理される。`do_pipe` の実際の作業はパイプにノードを割り当て、そのファイル記述子を 2 つ返すだけである。

バイブの読み書きも、ファイルの読み書きとは若干異なっている。バイブの容量は有限だからである。すでに満杯のバイブに書き込みを行おうとすると、保留状態になってしまう。同様に空のものもバイブから読出しも、保留状態を引き起こす。そのためバイブには、現在位置(読み出すとするものが使用)上、サイズ(書き込みとするものが使用)の2つのポインタが存在し、データをどこから書き込み、どこから読み出すかを決定できるようにになっている。

パイプに対する操作が可能かどうかを調べる各種検査が `pipe_check(10433 行目)` によって提供されている。呼出し側を保留状態にする可能性のあるこれらの検査以外にも、`pipe_check` を用いて `release` を呼び出し、データ型またはデータ不足によって保留になっていたプロセスが再開できるか調べることもできる。10457 行目と 10458 行目では、それぞれスレッド中の進出しおよび

- 指定した特殊ファイルはブロック型デバイスでない。
- 特殊ファイルはブロック型デバイスであるが、すでにマウントされている。
- ①マウントしようとしているファイルシステムに不正なマジックナンバーが与えられている。
- ②マウントしようとしているファイルシステムは無効である(例えばiノードがない)。
- ③マウントされる図のファイルが存在しない、または特殊ファイルである。
- マウントされたファイルシステムにビットマップのための空間がない。
- マウントされたファイルシステムにスーパーブロックのための空間がない。
- マウントされたファイルシステムにルートiノードのための空間がない。

おそらくこの様な点に固執することはあまり適当ではないように思えるだろう。しかし実践的なオペレーティング・システムでは必ずといっていいほど、かなりの量のコードをそれほど面白くも、また重要とも思われぬような作業に費やしているのが現状である。しかしこの様な作業こそシステムを向上させるのに欠かせないものなのである。ユーザーが1箇月に一度の割合で間違ったフロッピー・ディスクを誤ってマウントしようとし、これがクラッシュやファイルシステムの破壊を招くとしたら、ユーザーはこのシステムの信頼性に疑問を感じ、自分自身ではなくシステム設計者を非難するだろう。

トーマス・エジソンはこれに関連した言葉を残している。彼は「天才」が1%のひらめきと、99%の努力で作られるものだと言っている。良いシステムと、悪いシステムの違いはスケジューリング・アルゴリズムの差でなく、詳細部分への心配りにあるのである。

ファイルのマウント解除は、ファイルのマウントより簡単である。間違いない部分が少ないからである。重要な点は、削除されるファイルシステム上に、オープンしたままのファイルやワーキングディレクトリを持つプロセスが存在しないかを確認することである。この検査はいたって簡単である。iノードテーブル全体を走査し、メモリ内のiノードのうち削除しようとするファイルシステムに属しているものがないかを調べるだけでよい(ルートiノードを除く)。その様なiノードが存在すれば、**UMOUNT** システムコールは失敗に終わる。

mount.c の最後の行は、name_to_dev(11180 行目)である。この手続きは特殊ファイルを指定し、そのiノードを得て、その主デバイスおよび副デバイス番号を取り出す。これらはiノード内にあり、通常最初のゾーンが向かう場所に保存されている。このスロットは特殊ファイルがゾーンを持たないため、提供されている。

■ ファイルのリンクとリンク解除

次のファイルはファイルのリンクとリンク解除を行う link.c である。手続き do_link(11275 行目)は、大半のコードがエラーの検査用に費やされているという点で do_mount に似ている。

```
link(file_name, link_name);
```

上記の呼出しにおいて発生する可能性があるエラーの主なものを以下に挙げる。

- file_name が存在しない、またはアクセスできない。
- file_name がすでに最大数のリンクを行っている。
- file_name がディレクトリである(スーパーユーザーのみがリンクできる)。
- file_name がすでに存在する。
- file_name と link_name が異なるデバイスである。

エラーが存在しなければ、新しいディレクトリ・エントリが文字列 link_name と、file_name のiノード番号を使って作成される。実際の入力値は、searchdir によって 11324 行目の do_link から行われる。

ファイルはリンク解除することによって削除される。作業は do_unlink(11342 行目)によって行われる。ここでも各種の検査が最初に行われる。エラーがなければ、ディレクトリ・エントリがクリアされ、iノード内のリンク数が1つ減少される。

リンクカウンタがここで0になった場合、すべてのゾーンが truncate(11388 行目)によって解放される。この手続きはiノードから検出したゾーンそれぞれを解放するという、簡単な方法で処理されている。

5.7.6 その他のシステムコール

最後に、状態、ディレクトリ、保護、時間、その他各種の機能に関わる様々なシステムコールを見ていくことにする。

■ ディレクトリとファイル状態の変更

ファイル statdir.c には CHDIR、CHROOT、STAT、FSTAT という4つのシステムコールに関する処理が記述されている。10722 行目の last_dir では、パスの最初の文字が、スラッシュであるかを調べることによってパス検索が始まることを学んだ。結果によってポインタがワーキングディレクトリか、ルートディレクトリに設定される。

1つのワーキング(またはルート)ディレクトリから別なディレクトリに移ることは、呼出し順のプロセステーブルのこの2つのポインタを変更することである。この変更は do_chdir(11475 行目)と do_chroot(11500 行目)によって行われる。どちらも必要な検査を行ってから、change(11515 行目)を呼び出し、古いディレクトリに置き換えられる新しいディレクトリをオープンする。

11483 から 11490 行目のコードは、ユーザープロセスが行う CHDIR では実行されない、これはメモリアネージャ用に特別に設けられているもので、EXEC 処理を目的としてユーザーディレクトリに移ることができる。ユーザーが例えば a.out などのファイルで、自身のワーキングディレクトリ内で実行しようとする、メモリアネージャにとってはその場所を調べるよりも、そのディ

5章 ファイルシステム

レクトリに移動したほうが簡単である。

このファイル内で残っている残り2つのシステムコール(STAT)とSTATは、ファイルの指定方法を除けば、基本的には同じである。前者はパス名を使って指定し、後者はオープンされたファイルのファイル記述子を使用する。トップレベルの手続きである `do_stat` と `do_lstat` は、どちらも `stat_inode` を呼び出して作業を行う。 `stat_inode` を呼び出す前に、 `do_stat` はファイルを開くことで `i_inode` を指定する。この様にして `do_stat` と `do_lstat` は、 `stat_inode` に `i_inode` をポインタを渡す。

`stat_inode` の唯一の役割は `i_inode` から情報を取り出し、それをバッファ内にコピーすることである。 `stat_inode` は1つのメッセージ内に収められないため、 `rw_user` を使い、11624行目でコピーしなければならない。

■ 保護

MINIX の保護機構では、 `rw_x` ビットを使用している。この一連のビットは各ファイルに対して与えられており、それぞれ所有者と、そのグループ、そして一般ユーザーごとに分れている。ビットの値は `CHMOD` システムコールによって行う。このシステムコールは `do_chmod` (11677行) が実行する。一連の安全性検査を行った後、モードが11704行目で変更されている。

システムコール `CHOWN` は、特定のファイルの内部 `i_inode` (11715行目) はスーパーユーザーしか実行できないが、2つの実装方法も類似している。

システムコール `UMASK` を使って、ユーザーはマスク(プロセステーブル内に保存されている)を設定し、以降の `CREATE` システムコールにおいて保護ビットのマスクを行う。これは11752行目の `staterent` のみで行われる。ただしコードは古いマスク値の結果として返さなければならぬ。この様な余分な作業が、必要なコード数を3倍に増やしてしまう(11751から11753行目)。

システムコール `ACCESS` を用いて、プロセスは指定した方法で(例えば監視用)にファイルのアクセスを行えるかどうかを確認できる。 `do_access` (11760行目) を使って、ファイルの `i_inode` を得、内部手続き `forbidden` (11782行目) を呼び出して、アクセスが禁止されていないかを調べる。 `forbidden` は、 `uid` と `gid` のみではなく、 `i_inode` 内の情報も調べる。その内容によって `rw_x` グループのいずれかを選択し、アクセスの許可状態を調べる。

11782行目の `read_only` は、引数として渡される `i_inode` が存在するファイルシステムが、読み取り専用または読み書き両用のどちらかでマウントされているかを調べるものである。これは、読み取り専用でマウントされたファイルシステムに対する書き込みを防ぐために必要である。

5.7 MINIX ファイルシステムの表現

■ 時刻

MINIX には時刻に関するシステムコールがいくつか提供されている。 `UTIME`、 `TIME`、 `STIME`、 `TIMES` である。これらを図5.43にまとめておく。ファイルシステムが処理しなければならぬ理由は特になく、たまたまそうだったのである。

UTIME	ファイルの最終更新時刻をセット
TIME	現在の時刻を秒で調べる
STIME	時刻の設定
TIMES	プロセスのアクセスを調べる

図 5.43 時刻に関する4つのシステムコール

各ファイルに保持されているのは、ファイルが最後に修正された時刻を記録する32ビットの値である。この時刻は `i_inode` に保存されている。ファイルの所有者、またはスーパーユーザーは、システムコール `UTIME` を使いこれを更新できる。手続き `do_untime` (11877行目) は、 `i_inode` を得、ユーザーが指定した時刻をその `i_inode` に格納することによって、このシステムコールを実行している。

ファイルシステムでは時刻の保存を行っていない。カーネル内のクロックスタックがそれぞれを保存している。その結果、時刻を調べたり、設定するためにはクロックスタックにメッセージを送るしか方法はない。事実 `do_untime` と `do_stime` はどちらもそのために提供されているものである。時刻は1970年1月1日からの時間(エポックと呼ばれる)が、秒単位で保存されている。

アカウント情報もカーネル内に保存されている。各クロック刻みごとに、刻み分の課金(account)をあるプログラムに対して行う。この情報はシステムスタックにメッセージを送ることで設定できる。 `do_times` (11937行目) がこの作業を行っている。たいいていのCコンパイラは、外部シンドルの先頭に下線を付け、またたいいていのリンクではシンボルを8文字に切り捨てるため、 `do_untime` と `do_times` が識別できなくなってしまう。したがって、この手続きには、 `do_times` という名前を付けていない。

■ その他

ファイル `misc.c` には、他に適切な場所がないシステムコールの機能をいくつか記述している。システムコール `DUP` は、ファイル記述子を複製するものである。すなわちその引数と同じファイルを指す新しいファイル記述子を作成する。

`DUP` とは呼出し形式が若干異なるシステムコールもある。次の形式で呼出しを行う。

```
dup2(fd, fd2);
```

5 5 ファイルシステム

ここでfdはオープン済みのファイルを示すファイル記述子であり、fd2はファイルの割当てが行われていない整数である。この呼出しは、fd2をfdと同じファイルに対する有効ファイル記述子にするものである。

2つのシステムコールは同じメッセージ型を共有する。この2つは、DUP2のためにfdに設定されたO100ビットによって識別される。どちらのコールもdo_dup(11981行目)によって処理される。作業は簡単にファイル記述子とfileエントリを操作するだけである。

システムコールSYNCは、ディスクにロードされてから修正されたすべてのブロック、iノード、スーパーブロックをコピーするものである。処理はdo_sync(12018行目)によって行われる。該当するエントリを抽出するために、すべてのテーブルを検索する。

システムコールFORK、EXIT、SETは、実際にはメモリアネージャのシステムコールであるが、結果をここで提示しなくてはならない。プロセスがフォークを行うと、カーネル、メモリアネージャ、ファイルシステムに通知しなくてはならない。これらの「システムコール」は、ユーザープロセスからではなく、メモリアネージャから発行される。関連する情報を記録するのがその役目である。

システムコールではないが、miscの最後のdo_revive(12419行目)は、ファイルシステムが要求された作業(例えばユーザープロセスへの入力データ提供など)を行うことのできなかつたシステムタスクの作業を完了した時に呼び出される。ファイルシステムはここでプロセスを復元し、返答メッセージを送る。

5.7.7 入出力デバイスインターフェイス

MINIXにおける入出力操作は、カーネル内のタスクにメッセージを送信することによって行われている。これらのタスクに対するファイルシステムのインターフェイスは、ファイルdevice.c内に含まれている。このファイルには特殊ファイルに特別な操作を行う手続きも含まれている。最初に特殊ファイルをオープンすると、特別な操作を必要とする場合に備え、手続きdev_open(12233行目)が呼び出される。この手続きは特殊ファイルの主/副デバイス番号を取り出し、主デバイス番号をtable.cファイルのdmapテーブルに対するインデックスとして、ファイルシステム内の特別な処理を行うための手続きを呼び出す(12240行目)。通常、ここには特に何の作業も行わないno_callを設定するが、必要に応じてdmapに他の手続きを設定することもできる。

デバイスのクローズも同様である。この場合、dev_close(12261行目)が作業を行う。

実際にデバイスの入出力操作を必要とする場合は、dev_io(12261行目)が呼び出される。そして標準メッセージ(3章の図3.15参照)を作成し、それを指定したタスクに送る。read_writeがキャラクタ型特殊ファイルの処理用に呼び出され、rw_blockがブロック型特殊ファイル用に呼び出される。dev_ioがタスクからの回答を待っている間、ファイルシステムも待機する。内部にマルチプログラミング機能がないからである。ただしこの様な待機時間は通常非常に短いものである(最悪の場合でも2,300ミリ秒)。

5.7 MINIX ファイルシステムの實現

device.cでは、IOCTLと呼ばれるシステムコールだけを処理している。このシステムコールはタスク・インターフェイスと密接な関係にあるため、ここに配置されている。IOCTLを実行すると、do_ioctlが呼び出され、メッセージを作成した後、適切なタスクに送信する。

find_dev(12328行目)は、主/副デバイス番号をデバイス番号から抽出するためのものである。device.cの最後の3つの手続きは、ファイルシステム内では特に呼び出されることはない。3つともdmapから、間接的に呼び出される。読み書きにはrw_devまたはrw_dev2を使用する(コールの階層については図5.41を参照されたい)。何も行わない手続きが必要な場合は、no_callが呼び出される。

5.7.8 汎用ユーティリティ

ファイルシステムには各種の場所で使用されるいくつかの汎用ユーティリティが存在する。これらはファイルutility.c内にまとめられている。最初の手続きはclock_timeである。メッセージをクロックタスクに送り、現在の時刻を調べるものである。次の手続きはcmp_stringで、2つの文字列を比較する。そして次のcopyは、ファイルシステムのアドレス空間の一部のデータブロックを、別な場所に移動する。

手続きfetch_nameは、大半のシステムコールが引数としてファイル名を持っているために必要となる。ファイル名が与えられ、ユーザーからファイルシステムに対して送られるメッセージに含まれることができる。ファイル名が与えられ、ユーザー空間内にあるファイル名を指すポインタがメッセージに含まれる。fetch_nameはどちらの場合も検査を行い、結果として名前を得る。

no_syncはエラーハンドラであり、ファイルシステムに存在しないようなシステムコールを受け取った場合に呼び出される。最後にpanicは、メッセージを出力し、何か重大な問題が発生したことをカーネルに伝えるためのものである。

最後のファイルはputcである。2つの手続きがここに含まれているが、どちらもメッセージの印字に関わっている。標準のライブラリ関数は、ファイルシステムにメッセージを送信するため、使用することはできない。これら手続きはメッセージを直接端末タスクに送信する。

5.8 まとめ

ファイルシステムをその外観から眺めてみると、ファイルとディレクトリの集合、そしてそれを行うための操作から構成されているように見える。ファイルは読取り、書き込みが行え、ディレクトリは作成、削除が行える。そしてファイルを1つのディレクトリから別のディレクトリに移動することもできる。現在、多くのファイルシステムでは、階層型のディレクトリ・システムをサポートしており、ディレクトリが無限なくサブディレクトリを持つことができる。

内蔵から見たファイルシステムは、また違った顔を持っている。ファイルシステム設計者は記憶領域の割当て方法、ファイルとそれに含まれるブロックの管理などにも配慮しなければならぬ。またシステムによって各種のディレクトリ構造が存在することも学んだ。ファイルシステムの信頼性と性能も重要な問題であることがわかった。

最近のネットワーク・ファイルサーバー構築法についても、そのいくつかを学んだ。アトミック更新とトランザクションは、システムで最も重要な機能であることが多い。

セキュリティと保護は、システムユーザーにとっても、また設計者にとっても非常に重要な問題である。古くからあるシステムに見られるセキュリティの問題点と、システム全般に見られる一般的な問題点について解説した。またパスワード、アクセス制御リスト、その他の機能を用いた認証作業についても触れた。

最後に MINIX ファイルシステムの詳細について学んだ。そしてその規模は大きい、決して狭くはないことを知った。ユーザー・アクセスからの作業要求を受け入れ、手続きへのポインタのテーブルを使って、要求されたシステムコールを実行するための手続きを呼び出す。そのモジュール化された構造と、カーネルの外という位置付けによって、MINIX から切り離して若干の修正を加えれば、独立したネットワーク・ファイルサーバーとして用いることもできる。

練習問題

- あるオペレーティング・システムでは、システム全体で1つのディレクトリしかサポートしていないが、そのディレクトリが任意の長さのファイル名と、任意の数のファイルを持つことはできる。階層型のファイルシステムと呼べるようなものをシミュレートすることができるか、できるならば、その方法を示せ。
- ファイル/etc/passwd に対し、5つの異なるパス名を挙げてみよう(ヒント: ディレクトリ・エントリ、`..` と `...` を使用する)。
- 空ディレクトリ領域は空リストまたはビットマップを使って管理することができる。ディスクアドレスはDビットを必要としている。Bブロックから成るディスクのうち、Fブロックは空状態である。ここで空リストがビットマップより少ない空間を使用するような状態をあげてみよう。16ビット値をDに割り当てた場合、ディスクの利用率を計算せよ。
- MS-DOS コンピュータはFATを使ってディスクブロックの記録をとっている。その性能を、長いファイルに対するランダムなシークに関してUNIXと比較せよ。以下の2つの場合について答えよ。
 - (a) FATが必ずディスク上に存在する。
 - (b) FATが必ずメモリ内に存在する。
- UNIX でパス games/zapper をオープンするのに必要なディスク参照回数は?
- UNIX のiノードにおけるリンク数は冗長機能である。iノードを指しているディレクトリ・エントリの数を知らせるのが唯一の役割であるが、ディレクトリを調べれば同じことがわかる。それではなぜこの様な機能が設けられているのか。
- ある UNIX ファイルシステムでは、1024 バイトブロックと、16 ビットのディスクアドレスが用いられている。iノードはデータブロック用の8つのディスクアドレスと、1つの単一間接ブロックアドレス、そして1つの二重間接ブロックアドレスを格納できる。最大ファイルサイズを良く考えて求めよ。
- UNIX では2つのディレクトリが同じiノードを指すことにより、ファイルの共有が行われている。MS-DOS ではiノードが存在しない。MS-DOS では同時に2つのディレクトリにファイルをリンクすることができるか、できるならばその方法を説明せよ。

**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☒ **BLACK BORDERS**
- ☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☐ **FADED TEXT OR DRAWING**
- ☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKEWED/SLANTED IMAGES**
- ☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER:** _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.

THIS PAGE BLANK (USPTO)